

# Zero-Knowledge Proofs of Training for Deep Neural Networks

Kasra Abbaszadeh  
University of Maryland  
kasraz@umd.edu

Dimitrios Papadopoulos  
Hong Kong University of Science and Technology  
dipapado@cse.ust.hk

Christodoulos Pappas  
Hong Kong University of Science and Technology  
cpappas@connect.ust.hk

Jonathan Katz  
University of Maryland  
jkatz2@gmail.com

## ABSTRACT

A zero-knowledge proof of training (zkPoT) enables a party to prove that they have correctly trained a committed model based on a committed dataset without revealing any additional information about the model or the dataset. An ideal zkPoT should offer provable security and privacy guarantees, succinct proof size and verifier runtime, and practical prover efficiency. In this work, we present KAIZEN, a zkPoT targeted for deep neural networks (DNNs) that achieves the above ideals all at once. In particular, our construction enables a prover to iteratively train their model by the (mini-batch) gradient-descent algorithm where the number of iterations need not be fixed in advance; at the end of each iteration, the prover generates a commitment to the trained model attached with a succinct zkPoT, attesting to the correctness of the entire training process. The proof size and verifier time are independent of the iteration number.

KAIZEN relies on two essential building blocks to achieve both prover efficiency and verification succinctness. First, we construct an optimized GKR-style (sumcheck-based) proof system for the gradient-descent algorithm with concretely efficient prover cost; this scheme allows the prover to generate a proof for each iteration of the training process. Then, we recursively compose these proofs across multiple iterations to attain succinctness. As of independent interests, we propose a framework for recursive composition of GKR-style proofs and techniques such as aggregatable polynomial commitment schemes to minimize the recursion overhead.

Benchmarks indicate that KAIZEN can handle a large model of VGG-11 with 10 million parameters and batch size 16. The prover runtime is 22 minutes (per iteration), which is  $43\times$  faster than generic recursive proofs, while we further achieve at least  $224\times$  less prover memory overhead. Independent of the number of iterations and, hence, the size of the dataset, the proof size is 1.36 megabytes, and the verifier runtime is only 103 milliseconds.

## 1 INTRODUCTION

Machine learning with deep neural networks (DNNs) has received unprecedented attention in recent years. At the same time, this widespread use of neural networks has raised concerns about the provenance and integrity of DNN models; see, e.g., a recent blog post [1] from Google about integrity issues in the context of general machine learning models. Such concerns can arise at any stage of the model development process. In this work, we focus on the problems of *proofs of training* (PoTs) for ensuring the integrity of the training stage. Here, roughly speaking, a model owner wants to convince a verifier that a model  $\mathcal{M}$  was correctly trained—i.e., trained using a publicly known learning algorithm and public specifications of (e.g., the model architecture and batch size)—using

a specific dataset  $\mathcal{D}$ .<sup>1</sup> When the model and/or dataset are committed (i.e., not publicly available), one can also consider *zero-knowledge* PoTs (zkPoTs) that reveal no additional information about  $\mathcal{M}$  or  $\mathcal{D}$  beyond the correctness of the training procedure. At the most basic level, a (zk)PoT can be used by a model owner to substantiate the provenance of a model they release; this can be useful even when the model or the dataset is public since the cost to train a model can be orders of magnitude more expensive than verifying that training was correct. A (zk)PoT can also be useful for other applications:

- **Arbitrating copyright claims.** There is growing concern about models being trained on copyrighted data [35], and a model owner may want to prove that their model was trained *without* using some particular copyrighted data item. To achieve this goal, the model owner can commit to the dataset  $\mathcal{D}$  used to train the model, prove using known techniques [11, 47] that the data item in question is not present in  $\mathcal{D}$ , and then use a (zk)PoT to show that the model was trained using the committed dataset  $\mathcal{D}$ .
- **Distributed training.** A model owner might distribute the task of training a model across multiple untrusted workers [44]. In this setting, the model owner may send portions of the dataset to each worker and have them return partial results, which can then be aggregated into a final model. PoTs can allow workers to prove that the partial model they return was trained correctly.
- **Proof of ownership.** A deep learning model  $\mathcal{M}$  might be stolen or extracted using model-inference attacks [54]. In that case, the real model owner can prove that they were the ones who trained  $\mathcal{M}$  by committing to the dataset  $\mathcal{D}$  they used to train the model and then using a PoT to prove that  $\mathcal{M}$  was trained using  $\mathcal{D}$ . An adversary would not be able to generate such a proof without access to  $\mathcal{D}$ ; even with such access, they would have to invest the resources necessary to train the model from scratch.

An ideal zkPoT should have strong security guarantees, be succinct<sup>2</sup>, and not impose too much overhead on the prover. As we discuss next, however, known techniques fall short of achieving at least one of these properties when it comes to proofs of training for DNNs.

### 1.1 Prior Work and Limitations

Jia et al. [36] proposed a proof of training for DNNs. Unfortunately, their construction does not offer cryptographically strong security guarantees, and recent work [25, 26] has identified various concrete attacks on their scheme. Moreover, their construction is neither

<sup>1</sup>We note this is different from *proofs of inference* [27, 43, 45, 57] where, given a model  $\mathcal{M}$  and an input  $x$ , the goal is to prove that the classification of  $x$  was done using  $\mathcal{M}$ .

<sup>2</sup>Succinctness means the proof size and the verifier time are sublinear the computation size. For some applications it might be critical; consider e.g., a prover posting a (zk)PoT on blockchain [3], and it is desirable to minimize the on-chain verification overhead.

zero-knowledge nor succinct; it reveals model weights and data items to the verifier who partially re-executes training. Recently, Garg et al. [29] constructed zkPoTs with provable security; however, their construction is not succinct and, as a result, only supports basic learning algorithms such as logistic regression but cannot support DNNs due to the significant verification overhead.

In principle, one can build a zkPoT using generic zero-knowledge proofs. Unfortunately, despite recent advances in such constructions [6, 20, 34, 59], they are still not sufficiently scalable to practically instantiate zkPoTs for complex models like DNNs. In particular, succinct zero-knowledge proofs (zkSNARKs) [13, 38, 48] incur prohibitive prover costs in terms of both the proof-generation time (at least 1000× slowdown compared to the training time) and the required memory overhead. On the other hand, constructions of proofs with efficient provers [12, 30, 61] are not succinct.

As already noted, existing works on zero-knowledge proofs of inference [43, 45, 57] solve a different (and somewhat easier) problem than the one we are considering here. While some of the optimizations they propose are also useful in our setting, the techniques used in those works are not sufficient to yield efficient zkPoTs for DNNs—note, in particular, that DNN training involves multiple iterations, where each iteration can be roughly 100× more expensive than DNN inference. Eisenhofer et al. [23] used generic SNARKs to address the problem of verifiable machine unlearning, where the goal is to retrain a model by removing particular items from the dataset. The computational complexity of unlearning is fairly low; thus, generic proofs can efficiently support this problem.

## 1.2 Our Contributions

In this work, we propose a technique for constructing efficient and succinct zkPoTs for DNN models that are trained using multiple iterations of the mini-batch gradient-descent algorithm. We also implement and evaluate KAIZEN,<sup>3</sup> an instantiation of our approach. Conceptually, our work involves two high-level components:

**Efficient proofs of gradient descent.** First, we construct an optimized *proof of gradient descent* (PoGD), a scheme for proving the correctness of a single gradient-descent iteration. Our PoGD is constructed from GKR-style (i.e., sumcheck-based) proof systems [59, 62, 63] since they offer concretely efficient prover overhead and succinct verification. Moreover, such proofs can be optimized for matrix operations such as matrix multiplication or convolution used in gradient-descent for DNNs to obtain sublinear proof-generation time [45, 53]; this is in contrast to generic proofs, which have prover overhead at least linear in the complexity of the computation. More details about our PoGD are presented in Sections 3.1 and 4.

**Recursive composition of sumcheck-based proofs.** Although our PoGDs offer succinct proofs for each gradient-descent iteration, they do not suffice to obtain succinct proofs for DNN training overall because such training involves a number of iterations that are typically linear in the size of the dataset. Thus, having the PoT include a PoGD for each iteration would result in a linear-sized proof. Alternately, using a single GKR-style proof for all iterations would result in a linear prover memory and verification overheads; note such proofs have verification linear in the computation depth.

<sup>3</sup>KAIZEN is a Japanese term translated as “change for the better”.

Instead, we use recursive proof composition, also referred to as *incrementally verifiable computation* (IVC) [10, 14, 55], to achieve succinctness. Roughly, at each iteration we prove both that the iteration was performed correctly and (recursively) that there is a valid proof of correctness for all previous iterations. In this way, the proof generated at the final iteration demonstrates the correctness of the entire computation with a proof size and verification cost independent of the iteration number. Additionally, the incremental nature of an IVC allows for prover memory cost only proportional to the complexity of one iteration, not the entire computation.

The challenge in recursively composing our sumcheck-based PoGDs is the relatively large circuit size of the verifier algorithm, which imposes a significant recursion cost; this is in contrast to generic IVC schemes [14, 17, 39], which offer efficient recursion cost but incur substantial prover cost overall for the gradient-descent computation since they rely on generic zero-knowledge proofs. To overcome this, we propose a framework for minimizing the recursion overhead when building IVC from sumcheck-based proofs. In particular, inspired by Halo [17] and follow-up works [15], we use aggregation schemes to reduce the cost of verifying the polynomial commitments used in the PoGDs for all executed iterations.

A difficulty here is that sumcheck-based proofs use commitments to *multivariate* polynomials instead of the *univariate* polynomials involved in prior work [15, 17]; existing aggregation schemes for polynomial commitments only handle univariate polynomials, and naively extending them to support multivariate polynomials would result in quasilinear prover overhead. To mitigate this limitation, and as a key building block for our recursion framework, we propose an aggregation scheme for multivariate polynomial commitments with linear prover overhead and logarithmic verifier overhead. More details of our aggregation scheme are presented in Section 5.

We further improve the overhead of recursive composition of sumcheck-based proofs with sumcheck-specific optimizations for lightweight hash functions [7] and aggregation of the evaluations of circuit-wiring predicates; see Sections 3.2 and 6 for more details. To our knowledge, building IVC from sumcheck-based proofs has not been considered by prior work, and our framework and techniques might be of independent interest as stand-alone primitives.

**Implementation and evaluations.** We evaluate KAIZEN on several well-known architectures [41, 42, 51], e.g., a convolutional model VGG-11 with 10 million parameters, 11 layers, and batch size 16. The model is trained on the CIFAR-10 dataset. Our prover time is 22 minutes per iteration, and independent of the number of iterations and the dataset size, the proof size is 1.36 megabytes, and the verifier time is only 103 milliseconds. Moreover, we compare KAIZEN with generic IVCs [17, 21, 39] as the baseline, where we achieve 43× faster prover time and at least 224× less prover memory usage.

## 1.3 Organization of the Paper

In Section 2, we introduce notation and include some background information on zero-knowledge proofs and DNNs. In Section 3, we present an overview of our techniques. Section 4 describes our PoGD protocol. We introduce our aggregation scheme in Section 5, and use this as a key building block of our recursive sumcheck framework presented in Section 6. In Section 7, we present KAIZEN, which we evaluate and compare to prior generic IVCs in Section 8.

## 2 PRELIMINARIES

We use  $\mathbb{F}$  to denote a finite field. We let  $\lambda$  be the security parameter, and  $\text{negl}$  be a negligible function. We define  $[n] := \{0, 1, \dots, n-1\}$ . We use bold lowercase letters  $\mathbf{x}, \mathbf{y}$  for vectors, and bold uppercase letters  $\mathbf{X}, \mathbf{Y}$  for matrices. For a vector  $\mathbf{x}$ , we use both  $x_i$  and  $\mathbf{x}[i]$  to denote the  $i$ th element of  $\mathbf{x}$ , and for a matrix  $\mathbf{X}$ , both  $X_{i,j}$  and  $\mathbf{X}[i, j]$  indicate the element in row  $i$  and column  $j$ . We use  $\mathbf{X}[i, :]$  to denote the  $i$ th row and  $\mathbf{X}[:, j]$  to denote the  $j$ th column of the matrix  $\mathbf{X}$ . We use  $\circ$  for the Hadamard (element-wise) product.

**Merkle tree.** A Merkle tree [47] MT is a data structure that can be used to commit to a vector. It consists of the following algorithms:

- **Commit:** on input a vector  $\mathbf{x}$  and randomness  $r$ , returns the root  $\rho$  of the Merkle tree as a commitment to the vector.
- **Open:** on input a vector  $\mathbf{x}$ , randomness  $r$ , and opening index  $i$ , returns an element  $x_i$  and proof  $p_i$ .
- **Verify:** on input a commitment  $\rho$ , index  $i$ , element  $x_i$ , and proof  $p_i$ , returns accept or reject.

The commitment time is linear in the length of the vector; proofs and verifier times are logarithmic in the length of the vector.

**Multilinear extensions.** Let  $V : \{0, 1\}^\ell \rightarrow \mathbb{F}$  be a function. The *multilinear extension* of  $V$  is defined as the unique multilinear polynomial  $\tilde{V} : \mathbb{F}^\ell \rightarrow \mathbb{F}$  with  $\tilde{V}(\mathbf{x}) = V(\mathbf{x})$  for all  $\mathbf{x} \in \{0, 1\}^\ell$ . For  $\mathbf{b} \in \{0, 1\}^\ell$ , let  $\beta_{\mathbf{b}} : \{0, 1\}^\ell \rightarrow \{0, 1\}$  denote the function with  $\beta_{\mathbf{b}}(\mathbf{x}) = 1$  if  $\mathbf{b} = \mathbf{x}$  and  $\beta_{\mathbf{b}}(\mathbf{x}) = 0$  otherwise, and note that  $\tilde{\beta}_{\mathbf{b}}(\mathbf{x}) = \prod_{i=1}^\ell ((1 - x_i)(1 - b_i) + x_i b_i)$ . We can then compute  $\tilde{V}$  as

$$\tilde{V}(\mathbf{x}) = \sum_{\mathbf{b} \in \{0, 1\}^\ell} \tilde{\beta}_{\mathbf{b}}(\mathbf{x}) \cdot V(\mathbf{b}).$$

We can define multilinear extensions of vectors by viewing a vector  $\mathbf{v} = (v_0, v_1, \dots, v_{n-1}) \in \mathbb{F}^n$  as a function  $v : \{0, 1\}^{\log n} \rightarrow \mathbb{F}$  such that  $\forall i \in [n] : v(i) = v_i$ . A similar method is applied to the matrices.

### 2.1 Proofs, Arguments, and Commitments

An interactive proof for relation  $R$  with corresponding language  $L_R$  is an interactive protocol between a prover  $\mathcal{P}$  and a verifier  $\mathcal{V}$  on joint input  $x$ . The protocol satisfies *completeness* if, in an honest execution when the prover additionally holds a witness  $w$  with  $(x, w) \in R$ , the verifier always accepts. *Soundness* guarantees that a malicious  $\mathcal{P}$  cannot convince  $\mathcal{V}$  to accept when  $x \notin L_R$ , except with negligible probability. *Knowledge soundness* is a stronger property requiring that if a malicious prover can cause the verifier to accept with high probability on joint input  $x$ , then it is possible to extract a witness  $w$  from the prover such that  $(x, w) \in R$ . Such a protocol is called an *argument* if (knowledge) soundness holds against only computationally-bounded provers. Informally, a proof/argument is *zero-knowledge* if a malicious verifier learns no information from the protocol about  $w$  other than  $(x, w) \in R$ . We say a proof/argument is *succinct* if the runtime of  $\mathcal{V}$  and the communication between  $\mathcal{P}$  and  $\mathcal{V}$  are  $\text{poly}(\lambda, |x|, \log |w|)$ ; see Appendix A.1 for formal definitions.

A proof/argument is *public-coin* if the messages sent by the honest verifier consist simply of random challenges. Under certain conditions, a public-coin protocol can be made non-interactive using the Fiat-Shamir transform [28] in the random-oracle model by replacing the challenge for each round with the output of the random oracle evaluated on the transcript of prior rounds.

**Incrementally verifiable computation (IVC).** IVC schemes [55] enables succinct verification of iterative computations. Let  $\mathcal{F}$  be a function,  $\omega_0, \dots, \omega_{i-1}$  be witness auxiliary inputs, and  $z_0$  be an initial input; define  $\forall k \in [i] : z_{k+1} = \mathcal{F}(z_k, \omega_k)$ . IVC allows a prover to incrementally generate a proof  $\pi_i$  substantiating that there exist  $\omega_0, \dots, \omega_{i-1}$  such that  $z_i$  was correctly computed from  $z_0$ . More formally, an IVC protocol consists of the following algorithms:

- $\mathcal{G}$ : on input security parameter, returns public parameters pp.
- $\mathcal{P}$ : on input an iteration counter  $i$ , initial input  $z_0$ , last output  $z_{i-1}$ , auxiliary input  $\omega_{i-1}$ , proof  $\pi_{i-1}$ , and public parameters pp, returns the next iteration output  $z_i$  and a proof  $\pi_i$ .
- $\mathcal{V}$ : on input an iteration counter  $i$ , initial input  $z_0$ , last output  $z_i$ , proof  $\pi_i$ , and public parameters pp, returns accept or reject.

Each  $\pi_i$  is a zero-knowledge argument for the relation:

$$R_{\text{IVC}_i} = \left\{ \begin{array}{l} (i, z_0, z_i), (\omega_0, \dots, \omega_{i-1}) : \\ \forall k \in [i] : z_{k+1} = \mathcal{F}(z_k, \omega_k) \end{array} \right\}.$$

A canonical technique for constructing IVC is recursive composition of a baseline succinct non-interactive proof. In particular, let  $\mathcal{P}_b$  and  $\mathcal{V}_b$  be the prover and verifier algorithms of the baseline proof, respectively. Let  $\mathcal{F}_A$  be the augmented function defined as:

$$(\mathcal{F}(z_i, \omega_i), \mathcal{V}_b(i, z_0, z_i, \pi_i)) \leftarrow \mathcal{F}_A(i+1, z_0, z_i, \omega_i, \pi_i).$$

$\mathcal{P}$ , on input  $i+1, z_i, \omega_i$ , and  $\pi_i$ , outputs  $z_{i+1}$  and invokes  $\mathcal{P}_b$  to generate a proof  $\pi_{i+1}$  that  $\mathcal{F}_A(i+1, z_0, z_i, \omega_i, \pi_i) = (z_{i+1}, 1)$ . The incremental design of IVC ensures that prover memory overhead, proof size, and verifier overhead are independent of the number of iterations and only proportional to the complexity of  $\mathcal{F}_A$ .

**Polynomial commitments.** A polynomial commitment scheme (PCS) enables a prover  $\mathcal{P}$  to generate a commitment to a polynomial and later open the polynomial at some input evaluation point. More formally, a PCS consists of the following procedures:

- **KeyGen:** on input the security parameter, number of variables  $\ell$ , and degree bound  $d$ , returns public parameters pp.
- **Commit:** on input a polynomial  $f$ , randomness  $r$ , and public parameters pp, returns a commitment  $\sigma$  to the polynomial.
- **Open:** on input a point  $x$ , a polynomial  $f$ , the randomness  $r$ , and public parameters pp, returns the evaluation  $y = f(x)$  and an evaluation opening proof  $\pi$ .
- **Verify:** on input a commitment  $\sigma$ , point  $x$ , evaluation  $y$ , and evaluation opening proof  $\pi$ , returns accept or reject.

A PCS is *hiding* if the commitment reveals no information about the polynomial. A PCS satisfies *evaluation binding* if a prover cannot open a commitment  $\sigma$  to two distinct evaluations at any input point  $x$ . Knowledge soundness ensures that a prover generating a correct evaluation proof must know the polynomial underlying the commitment. A zero-knowledge PCS guarantees that evaluation proofs reveal no additional information about  $f$  other than the value  $f(x)$ . Formal definitions are provided in Appendix A.1.

### 2.2 GKR-Based Zero-Knowledge Arguments

Goldwasser et al. [31] proposed an interactive proof system for layered arithmetic circuits, referred to as the GKR protocol, based on the sumcheck protocol. We review the scheme here.

**The sumcheck protocol.** Fix an  $\ell$ -variate polynomial  $f : \mathbb{F}^\ell \rightarrow \mathbb{F}$  with variable-degree bound  $d$ . The sumcheck protocol [46] enables a verifier  $\mathcal{V}$  to delegate computation of the sum of  $f$  over the binary hypercube, i.e.,  $H = \sum_{b \in \{0,1\}^\ell} f(b)$ , to a prover  $\mathcal{P}$ . At the end of the protocol,  $\mathcal{V}$  needs the ability to learn the evaluation of  $f$  at a single (random) point  $r \in \mathbb{F}^\ell$ . The prover time is  $O(d^\ell)$ , and the verification time and proof size are  $O(d\ell)$ . The soundness error is  $O(d\ell/|\mathbb{F}|)$ . We present the scheme in Protocol 1 of Appendix A.

**The GKR protocol.** Let  $C$  be a circuit of depth  $d$  over a finite field  $\mathbb{F}$ , with layer 0 the output layer and layer  $d$  the input layer, and where each gate in the  $i$ th layer takes two inputs from the  $(i+1)$ st layer. Let  $S_i$  be the number of gates in the  $i$ th layer, and  $s_i := \lceil \log S_i \rceil$ . Define wiring predicates  $add_i, mult_i : \{0,1\}^{s_i+2s_{i+1}} \rightarrow \{0,1\}$ , where  $add_i$  (resp.,  $mult_i$ ) takes as input a gate label  $z$  for layer  $i$ , and gate labels  $x, y$  for layer  $i+1$ , then returns 1 iff gate  $z$  is an addition (resp., multiplication) gate with inputs  $x$  and  $y$ . Fixing some input to the circuit, we define  $v_i : \{0,1\}^{s_i} \rightarrow \mathbb{F}$  so that  $v_i(b)$  is the output of the  $b$ th gate in the  $i$ th layer. Note that

$$\begin{aligned} \tilde{v}_i(z) = & \sum_{x,y \in \{0,1\}^{s_{i+1}}} \widetilde{add}_i(z, x, y) \cdot (\tilde{v}_{i+1}(x) + \tilde{v}_{i+1}(y)) + \\ & \widetilde{mult}_i(z, x, y) \cdot \tilde{v}_{i+1}(x) \cdot \tilde{v}_{i+1}(y). \end{aligned} \quad (1)$$

$\mathcal{P}$  can prove the evaluation of  $C$  as follows. Assume  $\mathcal{V}$  has oracle access to  $\tilde{v}_0$  and  $\tilde{v}_d$  (i.e., the input and output).  $\mathcal{V}$  samples a random challenge  $r_0$  and evaluates  $\tilde{v}_0(r_0)$ . Then,  $\mathcal{P}, \mathcal{V}$  run the sumcheck protocol on Equation (1) for  $i = 0$ , at the end of which  $\mathcal{V}$  asks  $\mathcal{P}$  for two evaluations of  $\tilde{v}_1$ , say at points  $r_{1,0}, r_{1,1}$ . Say  $\mathcal{P}$  claims  $y_0 = \tilde{v}_1(r_{1,0})$  and  $y_1 = \tilde{v}_1(r_{1,1})$ . The verifier then chooses  $\gamma_0, \gamma_1 \leftarrow \mathbb{F}$  and the parties run the sumcheck protocol to prove that

$$\begin{aligned} \gamma_0 \cdot y_0 + \gamma_1 \cdot y_1 = & \sum_{x,y \in \{0,1\}^{s_{i+1}}} \\ & (\gamma_0 \cdot \widetilde{add}_i(r_{1,0}, x, y) + \gamma_1 \cdot \widetilde{add}_i(r_{1,1}, x, y)) \cdot (\tilde{v}_{i+1}(x) + \tilde{v}_{i+1}(y)) + \\ & (\gamma_0 \cdot \widetilde{mult}_i(r_{1,0}, x, y) + \gamma_1 \cdot \widetilde{mult}_i(r_{1,1}, x, y)) \cdot \tilde{v}_{i+1}(x) \cdot \tilde{v}_{i+1}(y). \end{aligned}$$

This reduces the problem to a claim of two evaluations of  $\tilde{v}_1$ . The parties then continue this procedure for  $d$  iterations; at the final iteration  $\mathcal{V}$  can check the evaluations of  $\tilde{v}_d$  using its oracle access. The formal protocol is provided in Protocol 2 of Appendix A. The GKR protocol can be made zero-knowledge by using zero-knowledge sumchecks [16, 59] and zero-knowledge polynomial commitments to provide oracle access to  $\tilde{v}_0$  and  $\tilde{v}_d$ . The protocol can be made non-interactive using the Fiat-Shamir transform.

## 2.3 Deep Neural Networks and Training

A DNN maps an input feature to a prediction by applying a sequence of transformations across  $L$  layers. The  $\ell$ th layer consists of a linear operation parameterized by weights  $W_\ell$  followed by application of a non-linear activation function. A layer is *dense* or *convolutional*. In dense layers, the linear operation returns  $T = W \cdot U$  given an input  $U$  and weights  $W$ , and convolutional layers returns  $T = W * U$ , where if  $U$  is of size  $u \times u$  and  $W$  is of size  $w \times w$ , then  $T$  is a  $(u - w + 1) \times (u - w + 1)$  matrix such that

$$T[i, j] = \sum_{a,b=0}^{w-1} U[i+a, j+b] \cdot W[a, b].$$

Following the linear operation, an activation function is performed; this is  $\text{ReLU}(x) = \max(x, 0)$ ,  $\tanh(x) = (e^{2x} - 1)/(e^{2x} + 1)$ , each applied coordinate-wise, or  $\text{Softmax}(x_i) = (e^{x_i})/(\sum_j e^{x_j})$ . A layer may also include a pooling that slides a *window* over the input and extracts one element from each window deterministically, e.g., it returns the maximum or average value in each of the windows.

**Training DNNs using mini-batch gradient-descent.** When training a DNN, the architecture, i.e., the number of the layers, the type of each layer, and an initial set of weights, is fixed in advance, and the goal is to optimize the weights at each layer to minimize a specified loss function  $\mathcal{L}$  over the training dataset. The training dataset is a set of data items, where each item consists of a feature and associated labels. The most popular optimization method for DNNs is mini-batch gradient-descent [50], which proceeds in a series of *epochs*. In each epoch, the dataset  $\mathcal{D}$  is randomly partitioned into batches  $B_1, \dots, B_k$ , each containing a specified number  $N$  of data items. An epoch consists of  $k$  iterations (one per batch). In the  $i$ th iteration, let  $U_{i,0}$  be the concatenation of the input features of the data items in  $B_i$  and let  $U_{i,\ell}$  (for  $\ell > 0$ ) be the concatenation of outputs of the  $\ell$ th layer given the current weights. In the  $i$ th iteration the current weights  $\{W_{i-1,\ell}\}_{\ell=1}^L$  are updated to obtain new weights  $\{W_{i,\ell}\}_{\ell=1}^L$  via the following steps:

- (1) *Forward pass:* We feed the inputs  $U_{i,0}$  to the model, apply layers sequentially, and for  $\ell = 1$  to  $L$ , compute  $U_{i,\ell}$ .
- (2) *Backward pass:* Let  $H_i$  be the average loss over the batch, and define the gradients  $G_{i,\ell} := \partial H_i / \partial W_{i-1,\ell}$  and  $R_{i,\ell} := \partial H_i / \partial U_{i,\ell-1}$ . These are computed for  $\ell = L$  to 1 as:

$$G_{i,\ell} = R_{i,\ell+1} \cdot \frac{\partial U_{i,\ell}}{\partial W_{i-1,\ell}}, R_{i,\ell} = \frac{\partial U_{i,\ell}}{\partial U_{i,\ell-1}} \cdot R_{i,\ell+1}. \quad (2)$$

$R_{i,L+1}$  is evaluated using the labels of the data items. Moreover, the gradients  $\partial U_{i,\ell} / \partial W_{i-1,\ell}$  and  $\partial U_{i,\ell} / \partial U_{i,\ell-1}$  are evaluated by applying some linear operations and non-linear activations to the inputs  $U_{i,\ell-1}$  and weights  $W_{i-1,\ell}$ , respectively.

- (3) *Update:* Given a learning rate  $\eta_\ell$  for each layer, the weights are updated as  $W_{i,\ell} = W_{i-1,\ell} - \eta_\ell \cdot \bar{G}_{i,\ell}$ , where  $\bar{G}_{i,\ell}$  indicates the average of  $G_{i,\ell}$  over the data items in the batch.

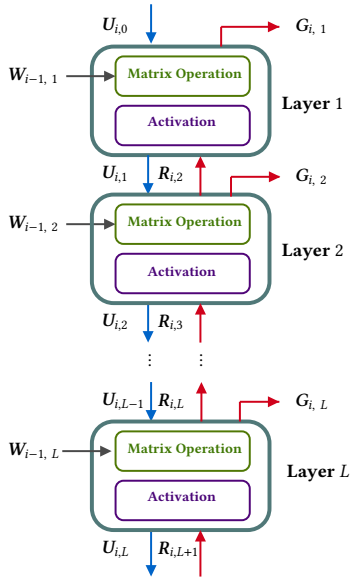
A high-level illustration is presented in Figure 1, and more details are presented in Algorithm 1 of Appendix A.

## 3 TECHNICAL OVERVIEW

In this section, we provide an overview of our techniques; low-level technical details are still deferred to subsequent sections.

### 3.1 Sumcheck Proofs for Gradient Descent

Our PoGD is a GKR-style sumcheck-based proof for one iteration of the mini-batch gradient-descent algorithm (cf. Section 2.3). In particular, given a fixed batch of data points and some model weights, PoGD enables the prover to update the weights by applying a gradient-descent iteration, and then, convince the verifier that the update is executed correctly. For the linear portions of the algorithm, we use optimized sumcheck-based proofs with sublinear proof-generation time. For the non-linear operations, we use bit decompositions along with a generic version of the GKR protocol. We discuss each of these components in more detail next.



**Figure 1: A high-level illustration of the  $i$ th gradient-descent iteration.**  $U_{i-1,0}$  denotes the input features,  $U_{i,\ell}$  (for  $\ell > 0$ ) denotes the output of layer  $\ell$ , and  $W_{i-1,\ell}$  denotes the current weights of layer  $\ell$ .  $R_{i,\ell}$  and  $G_{i,\ell}$  are the gradients computed in the backward pass. Blue and red arrows indicate the flow of the forward and backward passes, respectively.

**Linear operations.** Linear operations, such as matrix multiplication or convolution, are the most computationally intensive parts of the gradient-descent algorithm. Fortunately, prior work [45, 53] has shown optimized sumcheck protocols for these operations that have proof-generation time sublinear in the runtime of the computation itself. In particular, Thaler [53] proposed a sumcheck proof for the multiplication of two  $n \times n$  matrices with  $O(n^2)$  prover time; note the underlying computation takes time  $\Theta(n^3)$  if naive matrix multiplication is used, and even the best matrix multiplication algorithms require time  $\omega(n^2)$ . More recently, Liu et al. [45] gave a sumcheck protocol for convolution of an  $n \times n$  input matrix with a  $w \times w$  weights where the prover time is  $O(n^2 + w^2)$ , asymptotically faster than computing the convolution itself in  $O(n^2 \cdot w^2)$  time.

**Handling non-linear operations.** Sumcheck proofs only support arithmetic operations over finite fields. However, neural networks use non-linear operations that cannot be compactly encoded in an arithmetic circuit, even once the underlying values have been quantized (cf. Section 4.1). We use standard techniques to handle these. For example, proofs about a comparison such as  $a \geq 0$  can be handled using bit decomposition: the prover provides the binary representation of  $a$  (including a sign bit) as an auxiliary input, and the verification circuit checks that the binary representation is consistent. The sign bit is then taken as the result of the comparison. To reduce the size of the circuit and wiring predicates, we express consistency checks directly as sumcheck instances rather than reducing them to an arithmetic circuit. As another example, we handle exponentiations using piece-wise linear approximations, which have been shown to provide sufficient accuracy [19, 24, 49].

### 3.2 Recursive Composition of Sumcheck Proofs

A challenge in applying recursive proof composition to our PoGD is the relatively large size of the baseline verification circuit  $\mathcal{V}_b$  and hence  $\mathcal{F}_A$  (cf. the discussion of IVC in Section 2.1). The circuit  $\mathcal{V}_b$  mainly consists of three components: verifying opening proofs for polynomial commitments, verifying sumcheck messages, and verifying the hashes of messages used to generate the challenge for the Fiat-Shamir transform. We discuss how we improve the efficiency of each of these components.

**Commitment openings.** For a GKR-based zkSNARK, the verifier needs to verify openings of several polynomial commitments. It is thus desirable to have a polynomial commitment scheme with a small verification circuit. Known constructions fall short: schemes based on pairings [37] have  $O(1)$  verification overhead, but are not compatible with recursive proof composition [10], schemes based on Merkle trees [8, 32, 63] require the verifier to compute  $\Omega(\log^2 N)$  hash evaluations for a polynomial of size  $N$ , and constructions based on prime-order groups [18, 56] have verification time  $\Omega(\sqrt{N})$  making them impractical for recursive proof composition.

Several prior IVC schemes [15, 17] resolved this problem by employing commitment aggregation. Roughly speaking, rather than verifying the opening of commitments in each recursive step, the commitments across all steps are aggregated and the aggregate is verified at the end. Unfortunately, the aggregation scheme used in prior work only applies to polynomial commitment schemes for *univariate* polynomials, whereas for applications to GKR-based proofs we need commitments to *multivariate* polynomials. Naively extending existing univariate schemes to the multivariate case would incur superlinear, i.e.,  $\Theta(N \log N)$ , prover time.

In this work, we propose an aggregation scheme compatible with multivariate polynomial commitments. For polynomials with  $N$  coefficients, the scheme has  $O(N)$  prover time and  $O(\log N)$  proof size and verifier time. Our aggregation technique is inspired by the evaluation reduction technique [31], which is a specific case of aggregation where multiple evaluations of a polynomial are combined into a single evaluation. We generalize this idea to the case of having multiple committed polynomials with multiple evaluation points. Our techniques applies to various polynomial commitment schemes. We also make the proofs zero-knowledge.

**Sumcheck messages.**  $\mathcal{V}_b$  must evaluate the multilinear extensions of the wiring predicates at several random points. Evaluating these polynomials requires time linear in the size of the circuit. Interactive implementations of the GKR protocol mitigate this overhead by preprocessing the evaluations; the verifier samples the random points and computes the evaluations in advance [62, 63]. However, for recursive proof composition the baseline proof system must be non-interactive and such preprocessing is not possible.

We avoid this overhead using an approach similar to aggregation. The various wiring-predicate evaluations are combined into a single evaluation by a sumcheck run. The combined evaluation is passed to the next iteration, and verification is deferred until the end. While a similar effect could also be achieved using the evaluation reduction technique, which is a specific case of our commitment aggregation scheme, using the sumcheck protocol offers concretely better prover performance in this setting.

**Fiat-Shamir challenges.** Due to the Fiat-Shamir transform,  $\mathcal{V}_b$  must verify the evaluation of a hash function on many inputs. In particular,  $O(d_A \cdot \log |\mathcal{F}_A|)$  hashes need to be verified, where  $d_A$  denotes the depth of the augmented function  $\mathcal{F}_A$ . We use MiMC [4], a sumcheck-friendly hash, to reduce this cost. MiMC applies a low-degree round function for multiple rounds. As observed in prior work [7], we can represent each round of MiMC directly as a sumcheck instance, rather than by explicitly writing it as an arithmetic circuit, and thus further reduce the size of  $\mathcal{V}_b$ .

### 3.3 Constructing a zkPoT

We construct a zkPoT for DNNs by recursively composing our sumcheck-based PoGDs across multiple iterations. In doing so, there is one additional complication: we need to ensure that the prover uses a (random) partition of the committed dataset in each epoch. We ensure this as follows. The prover commits to the dataset using a Merkle tree. In each epoch, the prover then partitions the dataset using a public “quasi-random” permutation, e.g., a cubic function, applied to the indices of the  $n$  items in the dataset; the key for the permutation is varied for each epoch. During each iteration of an epoch, the prover provides the data items contained in the current partition along with their Merkle proofs. As part of computing gradient-descent iterations, the we additionally check the correctness of the partition as well as the Merkle proofs.

## 4 PROOFS OF GRADIENT DESCENT

We present our *proof of gradient descent* (PoGD), which is a GKR-style sumcheck-based proof for one mini-batch gradient descent iteration.

### 4.1 Handling Fixed-Point Operations

Model weights and data items are often represented by real numbers, but we need to encode them as field elements. We fix a field size  $p$  large enough to ensure no overflow during the computation. We quantize a real number  $r$  as a  $q$ -bit integer ( $q \ll p$ ) with  $f$ -bit precision as follows: when  $0 \leq r \leq 2^{q-f}$ , map  $r$  to  $s = \lfloor r \cdot 2^f \rfloor \in \mathbb{F}_p$ , and when  $-2^{q-f} \leq r < 0$ , map  $r$  to  $s = p - \lfloor r \cdot 2^f \rfloor \in \mathbb{F}_p$ .

**Integer arithmetic and comparison.** Two numbers can be added by adding their quantizations. For the multiplication of numbers  $a$  and  $b$  with quantizations  $s_a$  and  $s_b$ , respectively, we need to scale  $z = s_a s_b$  by a factor of  $2^{-f}$ . For this purpose, we ask the prover to provide  $z_0, \dots, z_{q-1}$ , the bit-decomposition of  $z$ , as an auxiliary input along with a sign bit  $z_q$  with  $z_q = 1$  indicating negativity; The verifier checks whether the provided bit-decomposition consists of binary values consistent with  $z$ , and returns the truncated output; i.e., if  $z = (-1)^{z_q} \sum_{i=0}^{q-1} z_i \cdot 2^i$  output  $(-1)^{z_q} \cdot \sum_{i=0}^{q-f-1} z_{i+f} \cdot 2^i$ .

A similar technique can be used for comparisons. Say we need to check whether  $z \geq 0$ . The prover can provide the bit-decomposition of  $z$ , including a sign bit; this can be easily checked by the verifier, and the result of the comparison is then the sign bit. We can also handle divisions. To compute the quotient  $a/b$ , where  $s_a$  and  $s_b$  are the quantizations of  $a$  and  $b$ , respectively, the prover provides a quotient  $z$  and a remainder  $r$ . The verifier then checks whether  $s_a = z \cdot s_b + r$  and  $r < s_b$  and returns  $z \cdot 2^f + r \cdot 2^{-f}$  as the output. To further reduce the circuit size, we accumulate scaling factors for each layer and defer scalings of to the output of the layer.

**Sumcheck for binary operations.** Scalings and bit-decomposition consistency checks can be handled by a field arithmetic circuit and verified by applying a generic proof to the circuit. However, for complex models, such binary operations blows up the size of the circuit and wiring predicates. To overcome this, we express these operations directly as sumcheck instances. In particular, let  $\mathbf{z}$  be an  $n$ -sized vector of values that the prover is required to provide their bit-decomposition. Let  $\mathbf{Z}_{\text{Bit}}$  be a matrix of size  $n \times q$  and  $\mathbf{z}_{\text{Sign}}$  be an  $n$ -sized vector, where  $\mathbf{Z}_{\text{Bit}}[i, \cdot]$  is claimed to be the bit-decomposition of  $\mathbf{z}[i]$ , and  $\mathbf{z}_{\text{Sign}}[i]$  indicates the associated sign bit. To check whether  $\mathbf{Z}_{\text{Bit}}$  and  $\mathbf{z}_{\text{Sign}}$  are binary, the verifier sends challenges  $r, r_x, r_y$ , and parties run the sumcheck protocol on:

$$0 = \sum_{x \in \{0,1\}^{\log n}, y \in \{0,1\}^{\log q}} \tilde{\beta}_x(r_x) \cdot \tilde{\beta}_y(r_y) \cdot (r \cdot \tilde{\mathbf{z}}_{\text{Sign}}(x) \cdot (1 - \tilde{\mathbf{z}}_{\text{Sign}}(x)) + \tilde{\mathbf{Z}}_{\text{Bit}}(x, y) \cdot (1 - \tilde{\mathbf{Z}}_{\text{Bit}}(x, y))). \quad (3)$$

The verifier then sends another challenge  $r_z$ , and parties verify the consistency of the bit-decompositions by running a sumcheck on:

$$\tilde{\mathbf{z}}(r_z) = \sum_{x \in \{0,1\}^{\log n}, y \in \{0,1\}^{\log q}} \tilde{\beta}_x(r_z) \cdot \tilde{\mathbf{Z}}_{\text{Bit}}(x, y) \cdot (1 - 2 \cdot \tilde{\mathbf{z}}_{\text{Sign}}(x)) \cdot 2^y. \quad (4)$$

We can handle scalings similarly. In particular, we only need to restrict the range of  $x$  and  $y$  to the proper intervals in Equation (4); the range of  $x$  selects which values are targeted and must be scaled, and the range of  $y$  determines what scaling factor we apply.

**Exponentiation.** Natural exponentiation is used in functions such as Softmax and tanh. In practice, weights and data items are typically normalized to the interval  $[-1, 1]$ . Similar to prior work [29, 49], we employ the following piece-wise approximation for  $e^x$ :

$$f(x) = \begin{cases} 0 & x < -1/2 \\ \frac{1}{2} + x & -1/2 \leq x \leq 1/2 \\ 1 & x > 1/2 \end{cases}$$

### 4.2 Our PoGD Design

The prover  $\mathcal{P}$  updates weights  $\mathbf{W}_{i-1}$  to  $\mathbf{W}_i$  using a data batch  $\mathbf{B}_{i-1}$  (cf. the discussion of gradient-descent in Section 2.3); the PoGD uses sumcheck proof messages to convince a verifier  $\mathcal{V}$  that the update is done correctly. For this purpose,  $\mathcal{P}$  uses a polynomial commitment scheme to give  $\mathcal{V}$  oracle access to the inputs and outputs of the algorithm (i.e.,  $\tilde{\mathbf{W}}_i, \tilde{\mathbf{W}}_{i-1}, \tilde{\mathbf{B}}_{i-1}$ ), and also auxiliary inputs denoted by  $\tilde{\mathbf{AUX}}_i$ , including bit-decompositions, quotients, and remainders; note that  $\mathbf{W}_i, \mathbf{W}_{i-1}$ , and  $\mathbf{AUX}_i$  are concatenations of  $\mathbf{W}_{i,\ell}, \mathbf{W}_{i-1,\ell}$ , and  $\mathbf{AUX}_{i,\ell}$ , respectively for all layers  $\ell = 1$  to  $L$ .

PoGD then proceeds in several phases. In the first phase, parties verify the update step of the algorithm, the second phase verifies the backward pass, and the third phase verifies the forward pass. Once sumcheck messages are generated, parties combine random evaluations of inputs and outputs received from different phases together by running additional sumcheck instances.  $\mathcal{V}$  verifies the final evaluations using its oracle access. For linear operations, parties run sublinear sumcheck protocols [45, 53], and for non-linear operation they run a generic GKR proof, e.g., Virgo++ [62] in our implementation. Below, we sketch the four phases of our PoGD.

**Phase 1:** For  $\ell = 1$  to  $L$ ,  $\mathcal{V}$  asks for a random evaluation of the trained weights  $\widetilde{W}_{i,\ell}$ . Having this evaluation, parties run GKR on  $W_{i,\ell} = W_{i-1,\ell} - \eta_\ell \cdot \widetilde{G}_{i,\ell}$ , where the run ends with random evaluations of the inputs  $\widetilde{W}_{i-1,\ell}$  and  $\widetilde{G}_{i,\ell}$ . The former proceeds to Phase 4, and the latter proceeds to Phase 1 to bootstrap the next sumchecks.

**Phase 2:** For  $\ell = 1$  to  $L$ , having the evaluation of  $\widetilde{G}_{i,\ell}$  received from Phase 1 and also the evaluation of  $\widetilde{R}_{i,\ell}$  from the previous iteration (If  $\ell \neq 1$ ) parties verify Equation (2). We use the sublinear sumcheck protocols for linear operations, GKR exploiting bit-decomposition of inputs for non-linear operations, and the sumcheck protocol on Equation (3) and Equation (4) for binary operations of scalings and consistency checks. The run ends with evaluations of  $\widetilde{W}_{i-1,\ell}$ ,  $\widetilde{U}_{i,\ell-1}$ ,  $\widetilde{R}_{i+1,\ell}$ , and  $\widetilde{\text{AUX}}_{i,\ell}$ , where  $\widetilde{R}_{i+1,\ell}$  proceeds to the next iteration,  $\widetilde{U}_{i,\ell-1}$  proceeds to Phase 3, and others proceed to Phase 4. For  $\ell = L$ , parties run GKR on the loss calculation from the actual labels of the data items included in the batch, yielding evaluations of  $\widetilde{B}_{i-1}$  and  $\widetilde{U}_{i,L}$  and proceeding to Phase 4 and Phase 3, respectively.

**Phase 3:** For  $\ell = L$  to 1, having evaluations of  $\widetilde{U}_{i,\ell}$  received from Phase 2 and the previous iteration (if  $\ell \neq L$ ), parties verify the forward pass by running the sublinear sumcheck protocols for the linear operations, GKR for non-linear operations, and the sumcheck protocol for scalings and consistency checks. The run ends with evaluations of inputs of the layer, i.e.,  $\widetilde{W}_{i-1,\ell}$ ,  $\widetilde{U}_{i,\ell-1}$ , and  $\widetilde{\text{AUX}}_{i,\ell}$ , where  $\widetilde{U}_{i,\ell-1}$  proceeds to the next iteration, and  $\widetilde{W}_{i-1,\ell}$  and  $\widetilde{\text{AUX}}_{i,\ell}$  proceed to Phase 4; note that  $\widetilde{U}_{i,0}$  is viewed as an evaluation  $\widetilde{B}_{i-1}$ .

**Phase 4:** Parties combine all evaluations of  $\widetilde{W}_{i-1}$ ,  $\widetilde{B}_{i-1}$ , and  $\widetilde{\text{AUX}}_i$  received from the previous phases. In particular, consider a generic  $\ell$ -variate multilinear polynomial  $f$  and multiple evaluations of form  $\{(x_i, y_i)\}_{i=1}^k$ , where  $f(x_i) = y_i$ . Given challenges  $\{\gamma_i\}_{i=1}^k$ , we reduce evaluations into a single one by running the sumcheck protocol on:

$$\sum_{i=1}^k \gamma_i \cdot y_i = \sum_{z \in \{0,1\}^\ell} f(z) \cdot \sum_{i=1}^k \gamma_i \cdot \widetilde{\beta}_z(x_i). \quad (5)$$

In Phase 4 of our PoGD, parties run this sumcheck by instantiating  $f$  with  $\widetilde{W}_{i-1}$ ,  $\widetilde{B}_{i-1}$ , and  $\widetilde{\text{AUX}}_i$ . At the end of the run, for each of the polynomials, the verifier receives a single evaluation. Along with the evaluation of  $\widetilde{W}_i$  from Phase 1, the verifier checks these using the oracle access, i.e., by asking for evaluation opening proofs.

We provide more formal details in Protocol 3 and present the security proof for our construction in Appendix B.

**Theorem 1.** *The PoGD construction presented in Protocol 3 satisfies completeness, knowledge soundness, and zero knowledge.*

Let  $s_{\text{in},\ell}$  be the input size and  $s_{\text{out},\ell}$  be the output size of the  $\ell$ th layer,  $N$  be the batch size,  $q$  be the quantization bit-length. We use Orion [60] as our polynomial commitment scheme due to its linear-time and concretely efficient prover and succinct verification. The prover time is  $O(\sum_{i=0}^{L-1} N s_{\text{in},\ell} s_{\text{out},\ell} + N q s_{\text{out},\ell})$ . The proof size and verifier time are  $O(\sum_{i=0}^{L-1} \log^2(N s_{\text{in},\ell} s_{\text{out},\ell}) + \log^2(N q s_{\text{out},\ell}))$ , which is dominated by evaluation opening proofs. In the next sections, we discuss how to reduce the verification cost using our aggregation scheme to  $O(\sum_{i=0}^{L-1} \log(N s_{\text{in},\ell} s_{\text{out},\ell}) + \log(N q s_{\text{out},\ell}))$  when we compose PoGDs recursively; see Sections 5 and 6. The protocol can be made non-interactive by the Fiat-Shamir heuristic.

## 5 AGGREGATABLE COMMITMENT SCHEMES

We present an aggregatable polynomial commitment scheme for *multivariate* polynomials, a key building block for our recursive proof composition framework. In our setting, a prover has generated commitments  $\sigma_1, \dots, \sigma_k$  to polynomials  $f_1, \dots, f_k : \mathbb{F}^\ell \rightarrow \mathbb{F}$  each is  $\ell$ -variate with variable-degree at most  $d$ . The prover then wants to convince a verifier that  $y_i = f_i(x_i)$  for some points  $x_1, \dots, x_k$  more efficiently than giving  $k$  independent evaluation opening proofs. We propose an scheme that allow the prover and verifier to map  $\{(\sigma_i, x_i, y_i)\}_{i=1}^k$  to a tuple  $(\sigma^*, x^*, y^*)$  such that  $\sigma^*$  is a commitment to a polynomial  $f^*$  of size at most as each input polynomial  $f_i$  and such that (with high probability)  $f^*(x^*) = y^*$  iff  $\forall i : f(x_i) = y_i$ ; thus a proof that  $f^*(x^*) = y^*$  can substantiate the original claim to the verifier. The stronger property of knowledge soundness ensures that a prover who can generate valid aggregation messages and an evaluation opening proof for  $(\sigma^*, x^*, y^*)$  must know  $f^*$  and the inputs  $\{f_i\}_{i=1}^k$ . The scheme is zero-knowledge if the process for generating  $(\sigma^*, x^*, y^*)$  and the eventual opening proof for  $f^*$  leak no additional information about the  $\{f_i\}_{i=1}^k$ . We present more formal definitions of these properties in Appendix C.

### 5.1 Aggregating Evaluations

We begin by reviewing a technique [31] for aggregating multiple evaluations of a single polynomial. We then show how to generalize that to support distinct polynomials, and to add zero-knowledge.

**Single polynomial, multiple points.** Consider the case where  $f_1 = \dots = f_k = f$  and  $\sigma_1 = \dots = \sigma_k = \sigma$ , but  $\{x_i\}$  are distinct. The prover and verifier interpolate a degree- $(k-1)$  univariate polynomial  $L : \mathbb{F} \rightarrow \mathbb{F}^\ell$  with  $L(i) = x_i$  for  $i \in \{1, \dots, k\} \subseteq \mathbb{F}$ . The prover then sends the verifier a polynomial  $g = f \circ L : \mathbb{F} \rightarrow \mathbb{F}$  of degree  $O(dk\ell)$ . The verifier checks that  $\forall i : g(i) = y_i$ , then, for a random  $r \in \mathbb{F}$ , the parties set  $x^* = L(r)$  and  $y^* = g(r)$ ; in this case, the output polynomial and commitments are  $f^* = f$  and  $\sigma^* = \sigma$ .

We claim that  $f(x^*) = y^*$  iff  $\forall i : f(x_i) = y_i$  with high probability. Completeness is immediate. Soundness can be argued as follows. Without loss of generality, say  $f(x_1) \neq y_1$ . If  $g = f \circ L$ , the verifier rejects since  $g(1) \neq y_1$ . If  $g \neq f \circ L$  then  $g$  and  $f \circ L$  agree on at most  $\deg(g) = O(dk\ell)$  points, and we have  $f(x^*) = y^*$  with a negligible probability at most  $O(dk\ell/|\mathbb{F}|)$ .

**Multiple polynomials, multiple points.** We now extend the above idea to the general case of (possibly) distinct polynomials  $\{f_i\}$  and points  $\{x_i\}$ . As before, the parties define  $L : \mathbb{F} \rightarrow \mathbb{F}^\ell$  such that  $\forall i : L(i) = x_i$ . The prover then sends the  $k$  univariate polynomials  $g_i = f_i \circ L$ , and the verifier checks whether  $\forall i : g_i(i) = y_i$ . Then, for random challenges  $\{\alpha_i\}_{i=1}^k$  and  $r$ , the parties set  $x^* = L(r)$ ,  $g = \sum_i \alpha_i g_i$ , and  $y^* = g(r)$ . The output polynomial is  $f^* = \sum_i \alpha_i f_i$ ; we defer to Section 5.2 the details of how the parties can generate the commitment  $\sigma^*$ ; we refer to that as *commitment aggregation*.

Again, we claim that  $f(x^*) = y^*$  iff  $\forall i : f(x_i) = y_i$  with high probability. Completeness holds because

$$f^*(x^*) = (f^* \circ L)(r) = \sum_i \alpha_i (f_i \circ L)(r) = \sum_i \alpha_i g_i(r) = g(r).$$

For soundness, if there exists some  $f_i(x_i) \neq y_i$ , the prover must send  $g_i \neq f_i \circ L$  to pass verification. But then with high probability over choice of  $\{\alpha_i\}, r$ , we have  $y = g(r) \neq (f^* \circ L)(r) = f^*(x^*)$ .



**Zero knowledge.** Before discussing commitment aggregation, we briefly discuss how to make the above idea zero-knowledge using masking polynomials. For all  $i$ , the prover randomly samples an  $\ell$ -variate polynomial  $h_i$  and sends a commitment  $\sigma'_i$  to  $h_i$  along with an evaluation  $v_i = h_i(x_i)$ . Given random challenges  $\{\beta_i\}$ , for all  $i$ , the prover then sends  $g_i = (f_i + \beta_i h_i) \circ L$ , and the verifier checks that  $g(i) = y_i + \beta_i v_i$ . The rest of the protocol is similar to above; here the output polynomial is  $f^* = \sum_i \alpha_i f_i + \alpha_i \beta_i h_i$ .

Naively, one could choose each  $h_i$ , above, as a random  $\ell$ -variate polynomial of variable-degree at most  $d$ . The number of coefficients in each  $h_i$  would then be  $O(d^\ell)$  and large, making it inefficient in practice to commit to the  $\{h_i\}$ . We show (cf. Appendix C) it suffices to sample each  $h_i$  as  $h_i = \sum_{j=1}^{\ell} h_{i,j}$ , where each  $h_{i,j}$  is a random univariate polynomial (in the  $j$ th variable) of degree  $d\ell$ . The size of each  $h_i$  is then  $O(d\ell^2)$ , and the cost of zero knowledge is  $O(kd\ell^2)$ .

## 5.2 Aggregating Commitments

What remains is to show how parties can compute  $\sigma^*$ , a commitment to  $f^* = \sum_i \alpha_i f_i$  given  $\{\sigma_i\}$ , commitments to the  $\{f_i\}$ ; the procedure naturally extends to the case of the zero-knowledge scheme, where  $f^* = \sum_i \alpha_i f_i + \alpha_i \beta_i h_i$ , and parties also hold  $\{\sigma'_i\}$ . Of course, if the underlying commitment scheme is homomorphic then commitment aggregation is trivial. However, practical homomorphic commitment schemes [18, 37, 56] are known only from group assumptions, and cannot be used efficiently for recursive proof composition. Indeed, in our application we use Orion [60], a commitment scheme based on Merkle trees that is not homomorphic. We describe how it is possible to commit to the linear combination of input polynomials for that scheme. Although the discussion is specific to Orion, the technique can be extended to other hash-based schemes [8, 63].

**Orion.** We briefly review Orion. In that scheme, the coefficients of a polynomial are mapped to the coefficients of a square matrix  $M$ , say of size  $n \times n$ . Then, a linear error-correcting code  $E: \mathbb{F}^n \rightarrow \mathbb{F}^p$  is applied to  $M$  twice. First, we apply the code to each row of  $M$ , yielding an  $n \times p$  matrix  $C_1$ . Then, we apply the code to each column of  $C_1$ , yielding a  $p \times p$  matrix  $C_2$ . The prover computes a Merkle root for each column of  $C_2$ , and then a Merkle root computed over those  $p$  roots is returned as a commitment to the polynomial.

The verifier checks the well-formedness of a commitment by sending a random vector  $r \in \mathbb{F}^n$  along with  $t = \Theta(\lambda)$  random indices in  $[p]$ . The prover responds with a vector  $v = r^T \cdot M$  along with the corresponding  $t$  columns of  $C_1$  and the associated Merkle proofs for the encodings of those columns. The verifier (1) evaluates the codeword  $w = E(v)$ ; (2) for each challenge index  $i$ , the returned column  $C_1[:, i]$  satisfies  $w_i = r^T \cdot C_1[:, i]$ ; and (3) the Merkle proofs are all correct. Evaluation proofs can be generated similarly, but we omit details. Orion achieves zero knowledge and succinctness by randomized encoding and composing the verification procedure described above with a generic succinct zero-knowledge proof [63].

**Aggregation for Orion.** Assume the verifier holds commitments  $\{\sigma_i\}$  to polynomials  $\{f_i\}$  known to the prover, and let  $M_i, C_{i,1}, C_{i,2}$  be the corresponding matrices that were generated by the prover when committing to  $f_i$ . The prover generates a commitment  $\sigma^*$  to  $f^* = \sum \alpha_i f_i$  and sends it to the verifier; let  $M^*, C_1^*, C_2^*$  be the corresponding matrices generated during this process. For an honest prover  $M^* = \sum_i \alpha_i M_i$ ,  $C_1^* = \sum_i \alpha_i C_{i,1}$ , and  $C_2^* = \sum_i \alpha_i C_{i,2}$ .

On the other hand, if  $M^* \neq \sum_i \alpha_i M_i$  then, letting  $\delta$  be the minimum distance of  $E$ , the matrices  $C_2^*$  and  $\sum_i \alpha_i C_{i,2}$  must differ in at least  $\delta^2$  entries.<sup>4</sup> The key observation is that, in addition to checking the well-formedness of  $\sigma^*$ , the verifier can check that the correct linear relationship holds between random entries of  $C_2^*$  and the corresponding entries of  $\{C_{i,2}\}$ ; if  $E$  has constant relative distance, it suffices to check that the linear combination relationship holds for  $\Theta(\lambda)$  randomly opened entries of each  $C_2^*$  and  $\{C_{i,2}\}$ .

We provide more formal details in Protocol 4 and prove the security of our aggregation scheme in Appendix C.

**Theorem 2.** *The aggregation scheme presented in Protocol 4 satisfies completeness, knowledge soundness, and zero knowledge.*

The prover cost of our scheme is dominated by generating the commitment  $\sigma^*$ , which can be done in  $O(d^\ell)$ , and evaluating the  $\{g_i\}$ , which can be done in  $O(k \cdot d^\ell)$  time. The proof size is dominated by the Merkle proofs, which have size  $O(kd\ell)$ ; this is also the verifier complexity. Moreover, the scheme can be made non-interactive via the Fiat-Shamir transform.

## 6 RECURSIVE SUMCHECK PROOFS

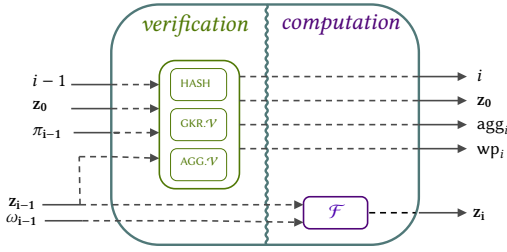
We present a generic framework for building IVC from GKR-style sumcheck-based proofs. In particular, our framework relies on the recursive proof composition technique when the baseline proof system  $(\mathcal{P}_b, \mathcal{V}_b)$  is instantiated by sumcheck-based non-interactive succinct proofs. When the baseline proof is recursively composed, the prover must generate a proof for the executions of the verifier algorithm  $\mathcal{V}_b$ . We discuss different components of this algorithm and our techniques to improve the efficiency of each components.

**Aggregation instead of opening.**  $\mathcal{V}_b$  verifies some sumcheck messages for the previous iteration of the augmented function  $\mathcal{F}_A$ . For this purpose,  $\mathcal{V}_b$  receives several random evaluations of the inputs and outputs of  $\mathcal{F}_A$ ; more precisely, evaluations of the multilinear extensions of the inputs and outputs. The inputs are committed by a succinct polynomial commitment scheme; such commitments prevent a blow-up in the input size across iterations. The prover then must provide opening proofs for the evaluations required by  $\mathcal{V}_b$  to verify sumcheck messages. Outputs of  $\mathcal{F}_A$  might be given to  $\mathcal{V}_b$  in plain or to be committed. Here, for simplicity, we assume they are given in plain, and they can be evaluated directly.

To reduce the complexity of  $\mathcal{V}_b$ , we replace opening proofs for committed polynomials with aggregation messages. In particular, for each  $i$ th iteration,  $\mathcal{V}_b$  receives two commitment/evaluation instances as follows. First,  $\mathcal{V}_b$  receives  $\text{agg}_{i-1} = (\sigma_{i-1}^*, x_{i-1}^*, y_{i-1}^*)$ , which is the aggregate instance returned the previous iteration; this includes a commitment  $\sigma_{i-1}^*$  and evaluation point  $(x_{i-1}^*, y_{i-1}^*)$ , and it is the aggregation of all commitments and their evaluations up to the iteration  $i - 2$ . Second,  $\mathcal{V}_b$  receives  $(\sigma_{i-1}, x_{i-1}, y_{i-1})$ , which is the instance for the inputs of the previous execution of  $\mathcal{F}_A$ . Here, again  $\sigma_{i-1}$  is a commitment to the inputs of  $\mathcal{F}_A$ , and  $(x_{i-1}, y_{i-1})$  is an evaluation point. These two instances are aggregated using our polynomial commitment aggregation scheme (cf. Section 5) and the aggregate output  $\text{agg}_i = (\sigma_i^*, x_i^*, y_i^*)$  is passed to the next iteration. In this way, we defer the opening to the final iteration.

<sup>4</sup>This assumes  $C_2^*, \{C_{i,2}\}$  were computed correctly, but in fact, it suffices to check that  $\sigma^*$  is well-formed to ensure that each  $\sigma_i$  were well-formed.





**Figure 2: High-level illustration of the augmented function  $\mathcal{F}_A$ .**  $\mathcal{F}$  denotes the iteration function. HASH, GKR.V, and AGG.V denote hash evaluation, sumcheck verifier, and aggregation verifier components of  $\mathcal{V}_b$ , respectively.  $\text{agg}_i$  denotes the  $i$ th aggregate commitment/evaluation instance and  $\text{wp}_i$  to denote the  $i$ th aggregate predicate evaluation.

**Handling wiring predicates.** In addition to the input and output evaluations,  $\mathcal{V}_b$  also requires several evaluations of the wiring predicates of  $\mathcal{F}_A$  at random points to validate sumcheck messages. Hard-coding all these predicates in  $\mathcal{V}_b$  blows up its circuit size and, hence, the recursion cost. Instead, we delegate these evaluations to the prover. We then apply an evaluation reduction sumcheck, i.e., the sumcheck protocol on Equation (5) to combine all wiring predicate evaluations into a single one, which is then passed to the next iteration. The final verifier receives the aggregation of wiring predicate evaluation for all iterations and validates its correctness.

**Sumcheck-friendly hashing.**  $\mathcal{V}_b$  verifies challenges for sumcheck proofs and aggregation messages as generated by the Fiat-Shamir transform. Furthermore,  $\mathcal{V}_b$  need to evaluate multiple hashes to verify Merkle proofs included in aggregation messages. To reduce the circuit size of this component, we instantiate the hash function with MiMC- $p/p$ , which is a sponge hash function, and each round function is  $F_i(x) = (x+k_i)^3 \bmod p$  for some round key  $k_i$ . Suppose that we have  $n$  number of hash evaluations and let  $d$  to denote the number of rounds. We can prove each hash round by running a single sumcheck protocol. Let  $v_0$  be the vector of hash outputs,  $v_d$  be the vector of inputs, and  $v_i := F_i(v_{i-1})$ . Given an initial random challenge  $r_0$ , the prover applies the sumcheck protocol on:

$$\tilde{v}_0(r_0) = \sum_{z \in \{0,1\}^{\log n}} \tilde{\beta}_z(r_0) \cdot (\tilde{v}_1(z) + k_1)^3, \quad (6)$$

The run ends with an evaluation  $\tilde{v}_1(r_1)$ , which proceeds the round. The prover continues this process to generate proof for all rounds.

**Putting everything together.** In the  $i$ th iteration,  $\mathcal{F}_A$  takes as input the iteration counter  $i$ , initial inputs  $z_0$ , last output  $z_{i-1}$ , an auxiliary input  $\omega_{i-1}$ , and proof  $\pi_{i-1}$ . The proof includes sumcheck proofs, aggregation messages, aggregate commitment/evaluation instance  $\text{agg}_{i-1}$ , aggregate wiring predicate evaluation  $\text{wp}_{i-1}$ , and a commitment and evaluation  $(\sigma_{i-1}, x_{i-1}, y_{i-1})$  from iteration  $i-1$ .  $\mathcal{F}_A$  invokes  $\mathcal{V}_b(i-1, z_0, z_{i-1}, \pi_{i-1})$ , which evaluates required hashes, verifies proof messages, aggregates  $\text{agg}_{i-1}$  with  $(\sigma_{i-1}, x_{i-1}, y_{i-1})$ , aggregates wiring predicate evaluations received at the end of the sumcheck verification with  $\text{wp}_{i-1}$ , and outputs the updated aggregates  $\text{wp}_i$  and  $\text{agg}_i$ . Once  $\mathcal{V}_b$  is completed successfully,  $\mathcal{F}_A$  then evaluates  $z_i \leftarrow \mathcal{F}(z_{i-1}, \omega_{i-1})$  and returns  $(i, z_0, \text{agg}_i, \text{wp}_i, z_i)$ .

We illustrate the augmented function in Figure 2. When execution of  $\mathcal{F}_A$  is completed, the IVC prover applies hash sumchecks, i.e., sumchecks on Equation (6), to MiMC- $p/p$  hashes and the baseline prover  $\mathcal{P}_b$  for other components. Furthermore, the prover needs to generate aggregation messages for the next iteration by running the prover of our aggregation scheme. In the final iteration, the prover additionally requires to generate an evaluation opening proof for the final aggregate commitment/evaluation instance.

In Protocol 5 of Appendix D, we provide more formal details of our recursive sumcheck framework and prove its security.

**Theorem 3.** *Protocol 4 presents an IVC construction that satisfies completeness, knowledge soundness, and zero knowledge.*

The complexity of the scheme depends on the baseline proof system and the polynomial commitment instantiation. For instance, consider an instantiation, where the GKR-style baseline proof is implemented by Virgo++ [62] using Orion [60] as the commitment; the prover then is  $O(|\mathcal{F}| + |z| + |\omega|)$ , and proof size and verifier are  $O(d \cdot \log |\mathcal{F}| + \log^2(|z| + |\omega|))$ , where  $d$  denotes the depth of  $\mathcal{F}$ .

## 7 ZERO-KNOWLEDGE PROOF OF TRAINING

We first define the notion of zkPoT as a cryptographic primitive and then propose KAIZEN a zkPoT for DNNs.

### 7.1 Definition

The prior work [29] proposed a monolithic definition of zkPoT, i.e., the prover executes a universal circuit, representing the entire training computation being verified, and generates a proof of the correct execution. Unfortunately, this approach is unsuited for iterative training algorithms as it requires fixing the number of training iterations in advance; the prover should complete the all iterations to generate a verifiable proof. For applications where a model owner incrementally updates their model, it is desirable to have a succinct zkPoT at any stage of the training process. To achieve this goal, we present an incremental definition of zkPoT, which generates a succinct zkPoT at the end of each iteration, attesting to the integrity of the entire training computation.

A consideration of such incremental definitions is ensuring data consistency; each iteration is only given a batch of data items claimed to be randomly sampled from the dataset. To resolve this, assume that there exists some fixed permutation of the dataset for each epoch and the dataset is committed to by a *position-binding* commitment, e.g., its Merkle root. At each iteration, the prover uses the batch consistent with the permutation and provides opening proofs to each data item in the batch. Before running the training iterations, Merkle openings are first validated. More precisely, let  $C$  be the iteration circuit in which the permutations are hard-coded. Let  $W_0$  be some initial model and  $\mathcal{D}$  be the dataset committed to by  $\rho$ . In each  $i$ th iteration,  $C$  takes as input the iteration counter  $i$ , dataset commitment  $\rho$ , the last model weights  $W_{i-1}$ , a batch of data samples  $B_{i-1} \subseteq \mathcal{D}$ , and the batch opening proof  $p_{B_{i-1}}$ .  $C$  first validates the opening proof and the consistency of data items in the batch with the hard-coded permutations. If verification passes, it applies the training algorithm and returns weights  $W_i$ . In practice, we can compactly hard-code permutations by fixing a key for an efficient pseudorandom permutation for each epoch, which enables  $C$  to evaluate the permuted positions required at each iteration.

A zkPoT enables the prover to first commit to the dataset and initial model weights, then sample a batch for each iteration and execute  $C$ . Once an iteration is completed, the prover can generate a commitment to the updated weights and a succinct proof that all iterations are correctly executed. More formally, a zkPoT consists:

- **KeyGen**: on input the security parameter, the setup algorithm generates the public parameters  $pp$ .
- **DataCom**: on input the training dataset  $\mathcal{D}$ , randomness  $r_{\rho_{\mathcal{D}}}$ , and public parameters  $pp$ , returns a commitment  $\rho_{\mathcal{D}}$ .
- **WeiCom**: on input some weights  $\mathbf{W}_i$ , a randomness  $r_{\sigma_{\mathbf{W}_i}}$ , and public parameters  $pp$ , returns a commitment  $\sigma_{\mathbf{W}_i}$ .
- **BatchOpen**: on input the dataset  $\mathcal{D}$ , randomness  $r_{\rho_{\mathcal{D}}}$ , the iteration counter  $i$ , and public parameters  $pp$ , returns the batch  $\mathbf{B}_{i-1}$  for the  $i$ th training iteration with an opening proof  $\mathbf{p}_{\mathbf{B}_{i-1}}$ .
- **Prove**: on input an iteration counter  $i$ , dataset commitment  $\rho_{\mathcal{D}}$ , initial weights commitment  $\sigma_{\mathbf{W}_0}$ , current weights  $\mathbf{W}_{i-1}$  with a commitment  $\sigma_{\mathbf{W}_{i-1}}$ , a batch  $\mathbf{B}_{i-1}$  with opening  $\mathbf{p}_{\mathbf{B}_{i-1}}$ , a proof  $\pi_{i-1}$ , and public parameters  $pp$ , returns weights  $\mathbf{W}_i$  and a proof  $\pi_i$ .
- **Verify**: on input an iteration counter  $i$ , the dataset commitment  $\rho_{\mathcal{D}}$ , the initial weights commitment  $\sigma_{\mathbf{W}_0}$ , the trained weights commitment  $\sigma_{\mathbf{W}_i}$ , and a proof  $\pi_i$ , returns accept or reject.

A proof-of-training satisfies soundness if the probability of passing verification for an adversary prover, returning inconsistent iteration counter and commitments is negligible. Knowledge soundness ensures that a prover accepted by the verifier must know satisfying weights and the dataset underlying the returned commitments. Zero knowledge guarantees that no information about the models and datasets is revealed; see Appendix E for formal definitions. For practicality, we require the proof size and the verifier runtime to be  $O(|C|)$  and independent of the number of iterations and the prover cost to be  $O(|C| \log |C|)$  per iteration.

## 7.2 Implementing a zkPoT for DNNs

KAIZEN commits to the dataset by its Merkle root, and the opening of a particular data batch includes Merkle proofs for each of the data items in the batch. Model weights are committed by generating a polynomial commitment to their multilinear extension.

Let PRP be a pseudorandom permutation, e.g., a cubic function, with an initial permutation key  $s_0$ . The key can be updated for each  $e$ th epoch to  $s_e = \mathcal{H}(e + s_0)$ , given a cryptographic hash function  $\mathcal{H}$ ; this ensures that each epoch trains the model on a randomly sampled partition. Let  $N$  be the batch size. For each  $i$ th iteration of  $e$ th epoch, the training iteration circuit  $C$  is executed as follows:

- (1)  $s_e \leftarrow \mathcal{H}(e + s_0), \forall k \in [N] : v_k \leftarrow \text{PRP}(s_e, (i-1)n + k)$ .
- (2) Validate  $\forall k \in [N] : \text{MT.Verify}(v_k, \rho_{\mathcal{D}}, \mathbf{B}_{i-1}[k], \mathbf{p}_{\mathbf{B}_i}[k]) = 1$ .
- (3) If the verification passes, return  $\mathbf{W}_i = \text{GD}(\mathbf{W}_{i-1}, \mathbf{B}_{i-1})$ , where GD denotes the the gradient descent algorithm (cf. Algorithm 1).

Consider an instantiation of our recursive sumcheck framework (i.e., Protocol 5), where the function  $\mathcal{F}$  is instantiated by  $C$  defined above. The baseline proof system is instantiated by the construction that verifies the execution of  $C$  by generating sumcheck messages for step 1 by a generic GKR-style proof [62], step 2 by running sumcheck for MiMC- $p/p$  hash (cf. Section 6), and step 3 via running our PoGD scheme as presented in Protocol 3.

A technical consideration is that in protocol 5,  $\mathcal{V}_b$  receives the initial inputs and the last iteration output in plain. In KAIZEN, these include model weights and can be significantly large. To prevent a blow-up in the size of  $\mathcal{V}_b$ , KAIZEN instead passes commitments to weights. The prover provides the evaluations of weights required for verifying sumcheck proofs, and these commitments and their evaluations are aggregated across the training iterations.

We present more formal details of the KAIZEN construction in Protocol 6 and a proof of its security in Appendix E.

**Theorem 4.** KAIZEN is a zero-knowledge proof of training and satisfies completeness, knowledge soundness, and zero knowledge.

Followed by the discussion in Section 4.2, let  $s_{\text{in},l}$  be the input and  $s_{\text{out},l}$  the output sizes of the  $l$ th layer,  $N$  the batch size, and  $q$  the bit-length of quantization. The KAIZEN prover cost per iteration is  $O(\sum_{l=0}^{L-1} N s_{\text{in},l} s_{\text{out},l} + N q s_{\text{out},l})$  both for generating PoGD messages and the recursion overhead. The final proof size and verifier cost are  $O(\sum_{l=0}^{L-1} \log^2(N s_{\text{in},l} s_{\text{out},l}) + \log^2(N q s_{\text{out},l}))$ , independent of the number of iterations and the size of the dataset.

## 8 IMPLEMENTATION AND EVALUATION

We discuss the implementation of KAIZEN and its evaluation results. The sumcheck and GKR protocols are implemented based on several open-source libraries [45, 62]. We run all experiments on a Linux VM with 8 physical CPUs, 2.80 GHz Intel(R) Xeon(R) Platinum 8370, and 256GB of RAM. Our current implementation is not parallelized, and we report the average running time of 10 executions for each of the benchmarks. Our code is publicly available [2].

**Protocol instantiation.** Our construction works on any finite field. We choose the extension field  $\mathbb{F}_{p^2}$  for the prime  $p = 2^{61} - 1$ . We use MiMC- $p^2/p^2$  [4] with 80 rounds and SHA-256 as our hash. The pseudorandom permutation PRP for shuffling the dataset is replaced by a cubic function over a field of the size as the dataset. We set  $q = 64$  and  $F = 32$  for the quantization scheme to ensure no overflow. The generic GKR-style proof used in our PoGD and KAIZEN protocols is instantiated by the Virgo++ [62]. Our parameters are chosen to achieve a  $\lambda = 100$ -bit security level.

Moreover, we instantiated the polynomial commitment scheme in our construction with Orion [60]. A difficulty with the Orion protocol is that its linear error-correcting code, i.e., Spielman [52], has a small relative distance, and this yields a significant number of Merkle tree openings for the aggregation scheme. In particular, given a relative distance  $\gamma = 0.15$ , as we use in our implementation,  $\lambda$ -bit security level requires 3045 openings. To reduce this number, instead of parsing the coefficients of polynomials as square matrices, we can parse them as matrices with a small constant number of columns and linear-sized rows. Although this yields a quasi-linear evaluation opening time, the commitment and aggregation times remain linear. As in KAIZEN, the evaluation opening is invoked only once at the final iteration, hence, the effect of this modification to the prover per iteration cost is negligible; here, we set the column-size to one to minimize the aggregation cost, yielding 426 openings.

Furthermore, we replace MiMC- $p^2/p^2$  with SHA-256 for up to 4 bottom layers of Merkle trees used in the Orion implementation to reduce the commitment-generation time for large polynomials.

**Gradient descent API.** We implemented an interference for the DNN gradient descent computation. A user can select specifications such as the model architectures, the batch size, and learning rate. The interference outputs sumcheck instances and GKR circuits, required for the mini-batch gradient-descent computation. The user can add any number of convolutional or dense layers to the architecture and specify input, output, and weight sizes. KAIZEN can support a wide range of activations such as ReLU, tanh, Sigmoid, and Softmax. For pooling methods, we can support both average and max pooling with arbitrary window sizes and strides.

**DNN Models and training datasets.** We measure the performance of our constructions on three well-known model architectures, LeNet [42], AlexNet [41], and VGG-11 [51]. For AlexNet and VGG-11, we truncated the number of parameters by reducing the number of channels and the size of the output size of the layers. LeNet has 3 convolutional layers, 2 average poolings, 2 dense layers, and a total number of 61706 trainable parameters. AlexNet has 5 convolutional layers, 3 average poolings, 2 dense layers, and a total number of 4.2 million parameters. Finally, the VGG-11 has 8 convolutional layers, 5 average pooling, 3 dense layers, and a total number of 10.1 million parameters. We use Softmax as the activation for the output layers and ReLU for others. We use the MNIST dataset [22], including 60,000 gray-scale  $32 \times 32$  images for LeNet, and CIFAR-10 [40], including 60,000 RGB  $64 \times 64$  images for AlexNet and VGG-11.

## 8.1 Performance of Our Constructions

We first measure the time the prover requires to generate sumcheck proofs of our PoGD at each iteration. Then, we report the overall prover and verifier overheads and the evaluation results of KAIZEN.

**Prover cost of PoGD sumchecks.** We report the prover time for sumchecks used in PoGD with a breakdown into the four phases of update step, backward pass, forward pass, and evaluation reduction. The batch size is ranged from  $N = 4$  to 16. The results are reported and summarized in Table 1; note that the prover cost of polynomial commitments are excluded, and we discuss them in the next reports. As shown in the table, it takes only 12.68 seconds to generate sumcheck messages for the gradient-descent of LeNet with  $N = 16$ . For AlexNet and VGG-11 it is 35.88 and 157.6 seconds, respectively.

**KAIZEN Performance.** Before starting the training procedure, the prover generates the Merkle root of the dataset. This takes only 18 seconds for MNIST and 218 seconds for CIFAR-10, where the Merkle tree is built using MiMC hash. Then, the prover executes the training iterations and generates proofs and commitments once an iteration is completed. We report the prover time per iteration with a breakdown into generating polynomial commitments, aggregation overhead<sup>5</sup>, generating sumcheck messages, and the overhead of proving the verifier function  $\mathcal{V}_b$ . Moreover, we report the proof size and verifier time with breakdown into sumcheck proofs and polynomial commitment evaluation openings proofs generated in the final iteration. We range  $N = 4$  to 16, and summarize in Table 2. As shown in the table, KAIZEN offers a prover runtime of 380 seconds per gradient-descent iteration for LeNet, 687 seconds for AlexNet, and 1357 seconds for VGG-11 given a batch size  $N = 16$ .

<sup>5</sup>The cost of generating commitments to aggregate polynomials returned by the aggregation scheme is considered part of the aggregation cost.

The verification cost is independent of the number of iterations and, hence, the size of the dataset. For batch size  $N = 16$ , The proof size is 900KB for LeNet, 1092KB for AlexNet, and 1360KB for VGG-11. The verifier time is only 66ms for LeNet, 81.84ms for AlexNet, and 103ms for VGG-11. This extremely fast verifier runtime is achieved by applying GKR-style proofs to the iteration circuit  $C$  and using Orion as our polynomial commitment scheme. As we observe, 67%-78% of the prover overhead is expended by the commitment and aggregation overheads. For some applications, it might be desirable to reduce the prover time at the cost of more verification time using different polynomial commitment schemes. For instance, one can instantiate KAIZEN with Bulletproofs [18] as the commitment scheme, which can achieve 2× faster prover at the cost of having verifier time in the order of hundreds of seconds.

## 8.2 Comparison to Generic IVCs

We compare the performance of KAIZEN with prior art generic IVC constructions, including Fractal [21], Halo [17], and Nova [39]. The generic constructions require the computation to be modeled as an arithmetic circuit. In particular, the schemes we discuss here rely on the *rank-1 constraint system* (R1CS) circuit representation. We carefully extracted the size of R1CS for the iteration circuit  $C$  to measure the performance of prior IVCs. For each operation of the gradient-descent algorithm, we consider the same computation procedure as our implementation, e.g., we use the standard matrix multiplication method, convolutions are turned into a few FFTs and evaluated by the Cooley–Tukey FFT algorithm, and non-linear operations are handled with bit-decompositions as auxiliary inputs. In addition to the circuit size, we also consider the output size of  $C$ , which includes model weights. Generic IVCs require the iteration outputs to be hashed by the verifier circuit multiple times. However, thanks to the use of GKR-style proofs, our construction can give the verifier only a commitment to the multilinear extension of outputs. We consider the hash function as Poseidon [33] for Nova and Fractal and Rescue [5] for Halo, as discussed in their protocol.

Given  $N = 16$ , the size of  $C$  for LetNet is about 250 million, for AlexNet is 618 billion, and for VGG-11 is 4.63 billion of gates; the weight hash constraint sizes are excluded. For this batch size, comparisons are reported in Table 3. As shown, for VGG-11, our prover is 43× faster than state-of-art generic IVC, which is Nova. Our verifier time is in the order of tens of milliseconds, comparable with Fractal and significantly faster than Halo and Nova. Our proof size is relatively large, however, for practical applications it might be considered negligible. Unfortunately, we are unable to run any of the mentioned IVC schemes even for LeNet with  $N = 1$ , as their memory usage exceeds our limitation, i.e., 256GB. In particular, we could only run Nova on circuit sizes up to 100 million gates, whereas LetNet for  $N = 1$  has 123 million gates, including 16 million gates for  $C$  and 107 million gates for weight hash. We extrapolated their performance based on the asymptotic complexities and their reported per-gate performance for small circuit sizes. Due to the substantial constants, reports are expected to be underestimated. Moreover, we measured the memory usage of the prior art (i.e., Nova) on circuit sizes up to 100 million gates. We observed that the memory usage of Nova grows linearly in the circuit size. Based on this, we expect to achieve 224× less memory overhead for VGG-11.

	LeNet			AlexNet			VGG-11		
	$N = 4$	$N = 8$	$N = 16$	$N = 4$	$N = 8$	$N = 16$	$N = 4$	$N = 8$	$N = 16$
Phase 1: Update	0.045	0.082	0.155	3.104	5.587	10.55	7.464	13.43	25.38
Phase 2: Backward	2.232	3.336	5.589	2.288	4.109	8.11	18.28	25.35	42.72
Phase 3: Forward	1.300	2.150	3.816	5.371	8.210	14.03	17.39	30.53	57.78
Phase 4: Evaluation	0.795	1.573	3.127	1.702	2.163	3.084	10.07	17.17	31.35
<b>Total Prover (s)</b>	<b>4.372</b>	<b>7.141</b>	<b>12.68</b>	<b>16.23</b>	<b>20.07</b>	<b>35.88</b>	<b>53.204</b>	<b>86.48</b>	<b>157.6</b>

Table 1: Prover runtime as a function of the batch size  $N$  when generating sumcheck messages in our PoGD.

	LeNet			AlexNet			VGG-11		
	$N = 4$	$N = 8$	$N = 16$	$N = 4$	$N = 8$	$N = 16$	$N = 4$	$N = 8$	$N = 16$
Commitments	73.37	97.87	146.7	84.5	140.7	253.1	141.4	254.4	480.3
Proof of Agg.	76.69	100.2	152.2	97.50	168.1	305.4	168.2	307.2	585.3
Proof of $V_b$	68.47	68.75	68.75	89.05	90.01	90.22	90.60	90.87	131.7
Proof of $C$	4.543	7.480	13.35	16.91	21.42	38.58	53.88	87.83	160.32
<b>Total Prover (s)</b>	<b>223.0</b>	<b>274.9</b>	<b>380.2</b>	<b>287.9</b>	<b>420.23</b>	<b>687.3</b>	<b>454.0</b>	<b>740.3</b>	<b>1357</b>
Sumcheck Proofs	457.7	472.7	487.7	619.1	637.1	656.0	837.6	861.7	886.8
Comm. Openings	366.8	434.2	412.2	388.9	395.9	436.3	424.4	448.3	473.3
<b>Total Proof (KB)</b>	<b>824.5</b>	<b>861.9</b>	<b>899.9</b>	<b>1008</b>	<b>1033</b>	<b>1092</b>	<b>1262</b>	<b>1310</b>	<b>1360</b>
Sumcheck Verifier	40.58	42.02	43.12	54.78	55.58	58.53	74.33	76.75	78.58
Comm. Verifier	19.92	20.81	22.85	21.00	22.14	23.31	22.13	23.31	24.84
<b>Total Verifier (ms)</b>	<b>60.50</b>	<b>63.01</b>	<b>65.97</b>	<b>75.78</b>	<b>77.72</b>	<b>81.84</b>	<b>96.46</b>	<b>100.06</b>	<b>103.12</b>

Table 2: The prover runtime per iteration, proof size, and verifier runtime of KAIZEN with batch size  $N$ .

	Prover (s)			Proof Size (KB)			Verifier (s)		
	LeNet	AlexNet	VGG-11	LeNet	AlexNet	VGG-11	LeNet	AlexNet	VGG-11
Fractal [21]	326,568	5,225,712	20,902,976	234	267	298	0.020	0.025	0.030
Halo [17]	50,850	813,595	3,254,382	4.86	5.09	5.45	3,970	63,528	254,114
Nova [39]	958	20,940	59,160	9.80	10.1	10.3	546	8,752	35,760
<b>KAIZEN</b>	<b>380</b>	<b>687</b>	<b>1,357</b>	<b>900</b>	<b>1092</b>	<b>1360</b>	<b>0.066</b>	<b>0.081</b>	<b>0.103</b>

Table 3: The prover runtime per iteration, proof size, and verifier runtime for KAIZEN and generic IVCs given batch size  $N = 16$ .

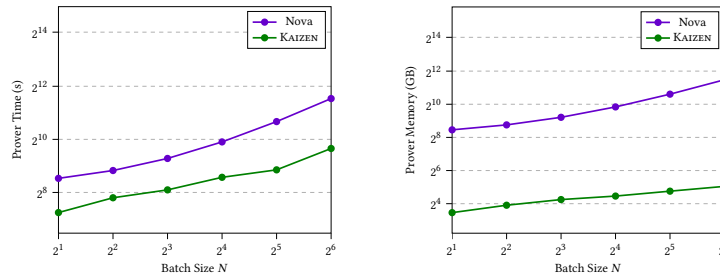


Figure 3: The prover time (per iteration) and memory overhead of KAIZEN and Nova for LeNet.

**More on KAIZEN vs. Nova.** We give a more detailed comparison between KAIZEN and Nova in terms of the prover runtime and memory usage. Nova offers better prover overhead compared to other generic IVCs. Specifically, we provide a comparison for the case of LeNet with batch sizes ranging from  $N = 2$  to 64. Results are reported in Figure 3. As shown, our prover is 2.4-3.7 $\times$  faster, and we require 31.9-86.42 $\times$  less memory overhead. As the batch size grows, the advantage of KAIZEN over Nova becomes more evident. We note that a larger batch size can also enable the prover to apply further optimizations by parallelizing both the training and proof generation computations in practice; e.g., prior work [58] showed that GKR-style proofs are amenable to high degrees of parallelism.

## ACKNOWLEDGMENTS

The work of Kasra Abbaszadeh and Jonathan Katz was supported by DARPA under Contract No. HR00112020025. The work Dimitrios Papadopoulos was supported Hong Kong RGC under grant 16200721. The views, opinions, and/or findings expressed are those of the authors and should not be interpreted as reflecting the position or policy of the Department of Defense or the U.S. Government, and no official endorsement should be inferred.

## REFERENCES

- [1] 2023. Increasing transparency in AI security. <https://security.googleblog.com/2023/10/increasing-transparency-in-ai-security.html>.

- [2] 2023. Kaizen. <https://github.com/zkPoTs/kaizen>.
- [3] 2023. Modulus Labs. <https://www.modulus.xyz>.
- [4] Martin Albrecht, Lorenzo Grassi, Christian Rechberger, Arnab Roy, and Tyge Tiessen. 2016. MiMC: Efficient Encryption and Cryptographic Hashing with Minimal Multiplicative Complexity. In *Advances in Cryptology—Asiacrypt 2016 (LNCS)*. Springer, 191–219.
- [5] Abdelrahman Aly, Tomer Ashur, Eli Ben-Sasson, Siemen Dhooghe, and Alan Szepieniec. 2020. Design of Symmetric-Key Primitives for Advanced Cryptographic Protocols. *IACR Trans. Symmetric Cryptol.* 2020 (2020), 1–45.
- [6] Scott Ames, Carmit Hazay, Yuval Ishai, and Muthuramakrishnan Venkatasubramanian. 2017. Liger: Lightweight Sublinear Arguments Without a Trusted Setup. In *ACM Conf. Computer and Communications Security 2017*. ACM, 2087–2104.
- [7] Alexandre Belling, Azam Soleimani, and Olivier Bégassat. 2023. Recursion over Public-Coin Interactive Proof Systems; Faster Hash Verification. In *ACM Conf. Computer and Communications Security*. ACM.
- [8] Eli Ben-Sasson, Iddo Bentov, Yinon Horesh, and Michael Riabzev. 2018. Fast Reed-Solomon Interactive Oracle Proofs of Proximity. In *Intl. Colloquium on Automata, Languages, and Programming (ICALP) (LIPIcs, Vol. 107)*. Schloss Dagstuhl—Leibniz-Zentrum für Informatik, 14:1–14:17.
- [9] Eli Ben-Sasson, Alessandro Chiesa, and Nicholas Spooner. 2016. Interactive Oracle Proofs. In *Theory of Cryptography (TCC)*. Springer, 31–60.
- [10] Eli Ben-Sasson, Alessandro Chiesa, Eran Tromer, and Madars Virza. 2014. Scalable Zero Knowledge via Cycles of Elliptic Curves. In *Advances in Cryptology—Crypto 2014*. Springer, 276–294.
- [11] Josh Benaloh and Michael de Mare. 1994. One-Way Accumulators: A Decentralized Alternative to Digital Signatures. In *Advances in Cryptology—Eurocrypt ’93*, Tor Helleseth (Ed.). Springer, 274–285.
- [12] Rishabh Bhadauria, Zhiyong Fang, Carmit Hazay, Muthuramakrishnan Venkatasubramanian, Tiancheng Xie, and Yupeng Zhang. 2020. Liger++: A New Optimized Sublinear IOP. In *ACM Conf. Computer and Communications Security*. ACM, 2025–2038. <https://doi.org/10.1145/3372297.3417893>
- [13] Nir Bitansky, Ran Canetti, Alessandro Chiesa, Shafi Goldwasser, Huijia Lin, Aviad Rubinfeld, and Eran Tromer. 2017. The Hunting of the SNARK. *J. Cryptol.* 30, 4 (oct 2017), 989–1066. <https://doi.org/10.1007/s00145-016-9241-9>
- [14] Nir Bitansky, Ran Canetti, Alessandro Chiesa, and Eran Tromer. 2013. Recursive Composition and Bootstrapping for SNARKs and Proof-Carrying Data. In *Symposium on Theory of Computing (STOC)*. ACM, 111–120.
- [15] Dan Boneh, Justin Drake, Ben Fisch, and Ariel Gabizon. 2021. Halo Infinite: Proof-Carrying Data from Additive Polynomial Commitments. In *Advances in Cryptology—Crypto 2021, Part I*. Springer, 649–680.
- [16] Jonathan Bootle, Alessandro Chiesa, and Katerina Sotiraki. 2021. Sumcheck Arguments and Their Applications. In *Advances in Cryptology—Crypto 2021, Part I*. Springer, 742–773. [https://doi.org/10.1007/978-3-030-84242-0\\_26](https://doi.org/10.1007/978-3-030-84242-0_26)
- [17] Sean Bowe, Jack Grigg, and Daira Hopwood. 2019. Recursive Proof Composition without a Trusted Setup. *Cryptology ePrint Archive*, Paper 2019/1021.
- [18] Benedikt Bünz, Jonathan Bootle, Dan Boneh, Andrew Poelstra, Pieter Wuille, and Greg Maxwell. 2018. Bulletproofs: Short Proofs for Confidential Transactions and More. In *IEEE Symposium on Security and Privacy*. IEEE, 315–334.
- [19] Gian Carlo Cardarilli, Marco Re, and Luca Di Nunzio. 2021. A pseudo-softmax function for hardware-based high speed image classification. *Sci. Reports* (2021).
- [20] Alessandro Chiesa, Yuncong Hu, Mary Maller, Pratyush Mishra, Noah Vesely, and Nicholas Ward. 2020. Marlin: Preprocessing zkSNARKs with Universal and Updatable SRS. In *Advances in Cryptology—Eurocrypt 2020*. Springer, 738–768.
- [21] Alessandro Chiesa, Dev Ojha, and Nicholas Spooner. 2020. Fractal: Post-quantum and Transparent Recursive Proofs from Holography. In *Advances in Cryptology—Eurocrypt 2020*. Springer, 769–793.
- [22] Li Deng. 2012. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine* 29, 6 (2012), 141–142.
- [23] Thorsten Eisenhofer, Doreen Riepel, Varun Chandrasekaran, Esha Ghosh, Olga Ohrimenko, and Nicolas Papernot. 2023. Verifiable and Provably Secure Machine Unlearning. [arXiv:2210.09126](https://arxiv.org/abs/2210.09126)
- [24] Alessandro Epasto, Mohammad Mahdian, Vahab Mirrokni, and Manolis Zampetakis. 2020. Optimal Approximation-Smoothness Tradeoffs for Soft-Max Functions. In *Neural Information Proc. Systems (NeurIPS)*. Curran Associates Inc.
- [25] Congyu Fang, Hengrui Jia, Anvith Thudi, Mohammad Yaghini, Christopher A. Choquette-Choo, Natalie Dullerud, Varun Chandrasekaran, and Nicolas Papernot. 2022. On the Fundamental Limits of Formally (Dis)Proving Robustness in Proof-of-Learning. <https://doi.org/10.48550/ARXIV.2208.03567>
- [26] Congyu Fang, Hengrui Jia, Anvith Thudi, Mohammad Yaghini, Christopher A. Choquette-Choo, Natalie Dullerud, Varun Chandrasekaran, and Nicolas Papernot. 2023. Proof-of-Learning is Currently More Broken Than You Think. In *European Symposium on Security and Privacy*. IEEE, 797–816.
- [27] Boyuan Feng, Lianke Qin, Zhenfei Zhang, Yufei Ding, and Shumo Chu. 2021. ZEN: An Optimizing Compiler for Verifiable, Zero-Knowledge Neural Network Inferences. *Cryptology ePrint Archive*, Paper 2021/087.
- [28] Amos Fiat and Adi Shamir. 1987. How To Prove Yourself: Practical Solutions to Identification and Signature Problems. In *Advances in Cryptology—Crypto ’86*, Andrew M. Odlyzko (Ed.). Springer, 186–194.
- [29] Sanjam Garg, Aarushi Goel, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoody, Guru-Vamsi Policharla, and Mingyuan Wang. 2023. Experimenting with Zero-Knowledge Proofs of Training. <http://ia.cr/2023/1345>
- [30] Irene Giacomelli, Jesper Madsen, and Claudio Orlandi. 2016. ZKBoo: Faster Zero-Knowledge for Boolean Circuits. In *Proceedings of the 25th USENIX Conference on Security Symposium (Austin, TX, USA) (SEC’16)*. USENIX Association, USA, 1069–1083.
- [31] Shafi Goldwasser, Yael Tauman Kalai, and Guy N. Rothblum. 2008. Delegating Computation: Interactive Proofs for Muggles. In *ACM Symposium on Theory of Computing*. ACM, 113–122.
- [32] Alexander Golovnev, Jonathan Lee, Srinath Setty, Justin Thaler, and Riad S. Wahby. 2021. Brakedown: Linear-Time and Post-Quantum SNARKs for R1CS. *Cryptology ePrint Archive*, Paper 2021/1043.
- [33] Lorenzo Grassi, Dmitry Khovratovich, Christian Rechberger, Arnab Roy, and Markus Schofnegger. 2021. Poseidon: A New Hash Function for Zero-Knowledge Proof Systems. In *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, 519–535.
- [34] Jens Groth. 2016. On the Size of Pairing-Based Non-interactive Arguments. In *Advances in Cryptology—Eurocrypt 2016*. Springer, 305–326.
- [35] Melissa Heikkilä. 2022. This Artist is Dominating AI-Generated Art. And He’s not Happy About It. MIT Technology Review.
- [36] Hengrui Jia, Mohammad Yaghini, Christopher A. Choquette-Choo, Natalie Dullerud, Anvith Thudi, Varun Chandrasekaran, and Nicolas Papernot. 2021. Proof-of-Learning: Definitions and Practice. In *IEEE Symposium on Security and Privacy*. IEEE, 1039–1056.
- [37] Aniket Kate, Gregory M. Zaverucha, and Ian Goldberg. 2010. Constant-Size Commitments to Polynomials and Their Applications. In *Advances in Cryptology—Asiacrypt 2010*, Masayuki Abe (Ed.). Springer, 177–194.
- [38] Joe Kilian. 1992. A Note on Efficient Zero-Knowledge Proofs and Arguments. In *Symposium on Theory of Computing (STOC)*. ACM, 723–732.
- [39] Abhiram Kothapalli, Srinath Setty, and Ioanna Tzialla. 2022. Nova: Recursive Zero-Knowledge Arguments from Folding Schemes. In *Advances in Cryptology—Crypto 2022*. Springer, 359–388.
- [40] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. [n. d.]. CIFAR-10 (Canadian Institute for Advanced Research). ([n. d.]). <http://www.cs.toronto.edu/~kriz/cifar.html>
- [41] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Adv. in Neural Information Processing Systems*, Vol. 25. Curran Associates, Inc.
- [42] Yan LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324. <https://doi.org/10.1109/5.726791>
- [43] Seunghwan Lee, Hankyung Ko, Jihye Kim, and Hyunok Oh. 2020. vCNN: Verifiable Convolutional Neural Network based on zk-SNARKS.
- [44] Mu Li, David G. Andersen, Jun Woo Park, Alexander J. Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J. Shekita, and Bor-Yiing Su. 2014. Scaling Distributed Machine Learning with the Parameter Server. In *Conference on Operating Systems Design and Implementation*. USENIX Association, 583–598.
- [45] Tianyi Liu, Xiang Xie, and Yupeng Zhang. 2021. zkCNN: Zero Knowledge Proofs for Convolutional Neural Network Predictions and Accuracy. In *Conf. on Computer and Communications Security*. ACM, 2968–2985.
- [46] Carsten Lund, Lance Fortnow, Howard Karloff, and Noam Nisan. 1992. Algebraic Methods for Interactive Proof Systems. *J. ACM* 39, 4 (oct 1992), 859–868.
- [47] Ralph C. Merkle. 1988. A Digital Signature Based on a Conventional Encryption Function. In *Advances in Cryptology—Crypto ’87*. Springer, 369–378.
- [48] Silvio Micali. 2000. Computationally Sound Proofs. *SIAM J. Computing* 30, 4 (2000), 1253–1298.
- [49] Payman Mohassel and Peter Rindal. 2018. ABY3: A Mixed Protocol Framework for Machine Learning. In *Conf. on Computer and Comm. Security*. ACM, 35–52.
- [50] Herbert Robbins and Sutton Monro. 1951. A Stochastic Approximation Method. *The Annals of Mathematical Statistics* 22, 3 (1951), 400 – 407.
- [51] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *Intl. Conf. on Learning Representations, ICLR*.
- [52] D.A. Spielman. 1996. Linear-time Encodable and Decodable Error-Correcting Codes. *IEEE Trans. Information Theory* 42, 6 (1996), 1723–1731.
- [53] Justin Thaler. 2013. Time-Optimal Interactive Proofs for Circuit Evaluation. In *Advances in Cryptology—Crypto 2013*. Springer, 71–89.
- [54] Florian Tramèr, Fan Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. 2016. Stealing Machine Learning Models via Prediction APIs. In *USENIX Security Symposium*. USENIX Association, USA, 601–618.
- [55] Paul Valiant. 2008. Incrementally Verifiable Computation or Proofs of Knowledge Imply Time/Space Efficiency. In *Theory of Cryptography (TCC)*. Springer, 1–18.
- [56] R. S. Wahby, I. Tzialla, A. Shelat, J. Thaler, and M. Walfish. 2018. Doubly-Efficient zkSNARKs Without Trusted Setup. In *IEEE Symposium on Security and Privacy*. IEEE, 926–943. <https://doi.org/10.1109/SP.2018.00060>

- [57] Chenkai Weng, Kang Yang, Xiang Xie, Jonathan Katz, and Xiao Wang. 2021. Mystique: Efficient Conversions for Zero-Knowledge Proofs with Applications to Machine Learning. In *USENIX Security Symposium*. USENIX Association, 501–518.
- [58] Tiancheng Xie, Jiaheng Zhang, Zerui Cheng, Fan Zhang, Yupeng Zhang, Yongzheng Jia, Dan Boneh, and Dawn Song. 2022. zkBridge: Trustless Cross-Chain Bridges Made Practical. In *Conf. on Computer and Communications Security*. ACM, 3003–3017.
- [59] Tiancheng Xie, Jiaheng Zhang, Yupeng Zhang, Charalampos Papamanthou, and Dawn Song. 2019. Libra: Succinct Zero-Knowledge Proofs with Optimal Prover Computation. In *Advances in Cryptology—Crypto 2019*. Springer, 733–764.
- [60] Tiancheng Xie, Yupeng Zhang, and Dawn Song. 2022. Orion: Zero Knowledge Proof with Linear Prover Time. In *Adv. in Cryptology—Crypto 2022*. Springer, 299–328.
- [61] Kang Yang, Pratik Sarkar, Chenkai Weng, and Xiao Wang. 2021. QuickSilver: Efficient and Affordable Zero-Knowledge Proofs for Circuits and Polynomials over Any Field. In *Conf. on Computer and Comm. Security*. ACM, 2986–3001.
- [62] Jiaheng Zhang, Tianyi Liu, Weijie Wang, YINUO Zhang, Dawn Song, Xiang Xie, and Yupeng Zhang. 2021. Doubly Efficient Interactive Proofs for General Arithmetic Circuits with Linear Prover Time. In *Conf. on Computer and Communications Security*. ACM, 159–177.
- [63] Jiaheng Zhang, Tiancheng Xie, Yupeng Zhang, and Dawn Song. 2020. Transparent Polynomial Delegation and Its Applications to Zero Knowledge Proof. In *IEEE Symposium on Security and Privacy*. 859–876.

## A ADDITIONAL PRELIMINARIES

In this section, we provide additional preliminaries, including formal definitions of arguments, polynomial commitment, protocols of sumcheck and GKR, and more of an in-detail presentation of the gradient descent algorithm for neural networks.

### A.1 Arguments and Commitments

**Argument systems.** An argument satisfies completeness if an honest prover always gets accepted. The scheme satisfies knowledge soundness if, for any accepted prover, there exists an extractor that returns the underlying witness having oracle access to the prover. Zero-knowledge ensures that a simulator can generate a transcript indistinguishable from a real protocol execution without having any knowledge of the witness. We use  $\mathcal{G}$  to represent the setup that generates public parameters pp.

**Definition 1. (Zero-knowledge arguments)** Given an NP relation  $R$  with language  $L_R$ , a zero-knowledge argument of knowledge for  $R$  is tuple of algorithms  $(\mathcal{G}, \mathcal{P}, \mathcal{V})$  if the following holds.

- **Completeness.** For every  $(x, w) \in R$  and  $\text{pp} \leftarrow \mathcal{G}(1^\lambda)$ ,

$$\Pr[\langle \mathcal{P}(\text{pp}, w), \mathcal{V}(\text{pp}) \rangle(x) = 1] = 1$$

- **Knowledge soundness.** For any PPT adversary  $\mathcal{A}$ , there exist a PPT extractor algorithm  $\xi^{\mathcal{A}}$  such that for any  $x$  and  $\text{pp} \leftarrow \mathcal{G}(1^\lambda)$ ,

$$\Pr[w \leftarrow \xi^{\mathcal{A}}(\text{pp}, x) : (x, w) \in R] \geq \Pr[\langle \mathcal{A}(\text{pp}), \mathcal{V}(\text{pp}) \rangle(x) = 1] - \text{negl}(\lambda)$$

- **Zero knowledge.** There exists a PPT simulator  $\mathcal{S}$  such that for every adversary  $\mathcal{A}$ ,  $(x, w) \in R$ , and  $\text{pp} \leftarrow \mathcal{G}(1^\lambda)$ ,

$$\text{View}(\langle \mathcal{P}(\text{pp}, w), \mathcal{A}(\text{pp}, z) \rangle(x)) \approx \mathcal{S}^{\mathcal{A}}(x, z)$$

An argument system is said to be transparent if the public parameter pp is simply a uniform random string, i.e., there is no secret state to generate pp.

**Polynomial commitment schemes.** The conventional definition of polynomial commitment promises binding and hiding security notions. Binding can be considered commitment binding, where a commitment cannot be opened at two polynomials, or evaluation binding, where a tuple of a commitment and an evaluation cannot be opened at two input points. Hiding ensures that the commitment reveals no information about the polynomial. However, we require stronger guarantees in this work. We rely on commitments that satisfy knowledge soundness, a stronger notion than binding, and zero-knowledge, which implies the hiding property.

Knowledge soundness ensures that an accepted prover must know the polynomial underlying the commitment and evaluation proofs. More precisely, for any accepted prover, there is an extractor that has oracle access to the prover and can extract the underlying polynomial. This notion implies both the commitment and the evaluation binding. Zero-knowledge guarantees that commitment and evaluation proofs reveal no additional information about the polynomial. Similar to arguments, we provide a simulation-based definition for zero-knowledge. Trivially, zero-knowledge implies hiding when the prover honestly generates the commitment.

### Definition 2. (Zero-knowledge polynomial commitment)

Given a finite field  $\mathbb{F}$  and an  $\ell$ -variate polynomial  $f$  with degree  $d$ , a zero-knowledge polynomial commitment scheme PCS is tuple of algorithms (KeyGen, Commit, Open, Verify) if the following holds.

- **Completeness.** For input  $x \in \mathbb{F}$  and public parameters pp,

$$\Pr \left[ b = 1 \mid \begin{array}{l} \sigma \leftarrow \text{Commit}(f, r, \text{pp}); \\ (y, \pi) \leftarrow \text{Open}(f, r, x, \text{pp}); \\ b \leftarrow \text{Verify}(\sigma, x, y, \pi, \text{pp}) \end{array} \right]$$

- **Knowledge soundness.** For any PPT adversary  $\mathcal{A}$ , there exist a PPT extractor algorithm  $\xi^{\mathcal{A}}$  such that given public parameters pp,

$$\Pr [ f \leftarrow \xi^{\mathcal{A}}(\text{pp}) : f(x) = y ] \geq \Pr \left[ \begin{array}{l} (\sigma, x, y, \pi) \leftarrow \mathcal{A}(\text{pp}); \\ \text{Verify}(\sigma, x, y, \pi, \text{pp}) = 1 \end{array} \right] - \text{negl}(\lambda)$$

- **Zero knowledge.** There exists a simulator  $\mathcal{S}$  such that for every adversary  $\mathcal{A}$  the two following experiments are indistinguishable.

$\text{Real}_{\mathcal{A}, f}(1^\lambda)$ (1) $\text{pp} \leftarrow \text{KeyGen}(1^\lambda, \ell, d)$ (2) $\sigma \leftarrow \text{Commit}(f, r, \text{pp})$ (3) $x \leftarrow \mathcal{A}(\sigma, \text{pp})$ (4) $(y, \pi) \leftarrow \text{Open}(f, r, x, \text{pp})$ (5) $b \leftarrow \mathcal{A}(\sigma, x, y, \pi, \text{pp})$ (6) Return $b$
$\text{Ideal}_{\mathcal{A}, \mathcal{S}}(1^\lambda)$ (1) $(\sigma, \text{pp}) \leftarrow \mathcal{S}(1^\lambda, \ell, d)$ (2) $x \leftarrow \mathcal{A}(\sigma, \text{pp})$ (3) $\pi \leftarrow \mathcal{S}(\sigma, x, f(x), \text{pp})$ (4) $b \leftarrow \mathcal{A}(\sigma, x, f(x), \pi, \text{pp})$ (5) Return $b$

In particular, for any  $\mathcal{A}$ , we have

$$|\Pr[\text{Real}_{\mathcal{A}, f}(1^\lambda) = 1] - \Pr[\text{Ideal}_{\mathcal{A}, \mathcal{S}}(1^\lambda) = 1]| \leq \text{negl}(\lambda)$$

### A.2 The Sumcheck and GKR Protocols

We present the sumcheck scheme in Protocol 1. Moreover, we provide details of the GKR interactive proof for layered circuits in Protocol 2. In our implementation, we use a variant of the GKR sumcheck-based proof that can support non-layered circuits with gates having arbitrary fan-in [62].

**Theorem 5.** [46] *Given a finite field  $\mathbb{F}$  and an  $\ell$ -variate polynomial  $f : \mathbb{F}^\ell \rightarrow \mathbb{F}$  with degree bound  $d$ , Protocol 1 is an interactive proof for the statement  $H = \sum_{b \in \{0,1\}^\ell} f(b)$  with soundness error  $O(\frac{d\ell}{|\mathbb{F}|})$ . It uses  $O(\ell)$  of interaction, running time of  $\mathcal{P}$  is  $O((d+1)^\ell)$ , and the proof size and the running time of  $\mathcal{V}$  is  $O(d\ell)$ .*

**Theorem 6.** [59] *Let  $C : \mathbb{F}^{\text{sin}} \rightarrow \mathbb{F}^{\text{sout}}$  be a depth- $d$  layered arithmetic circuit. Protocol 2 is an interactive proof for  $C$  with soundness error  $O(d \cdot \frac{\log |C|}{|\mathbb{F}|})$ . It uses  $O(d \cdot \log |C|)$  of interaction and running time of  $\mathcal{P}$  is  $O(|C|)$ . Let  $T$  be the optimal computation time for all  $\widetilde{\text{add}}_i$  and  $\widetilde{\text{mult}}_i$ . The running time of  $\mathcal{V}$  is  $O(s_{\text{in}} + s_{\text{out}} + d \cdot \log |C| + T)$ . For log-space uniform circuits, it is  $T = \text{polylog}|C|$ .*



### Protocol 1. Sumcheck

The protocol proceeds in  $\ell$  rounds.

(1) In the first round,  $\mathcal{P}$  sends the polynomial

$$f_1(x_1) := \sum_{b_2, \dots, b_\ell \in \{0,1\}} f(x_1, b_2, b_3, \dots, b_\ell)$$

$\mathcal{V}$  checks if  $H = f_1(0) + f_1(1)$  and sends  $r_1 \leftarrow \mathbb{F}$ .

(2) In the  $i$ th round ( $2 \leq i \leq \ell - 1$ ),  $\mathcal{P}$  sends the polynomial

$$f_i(x_i) := \sum_{b_{i+1}, \dots, b_\ell \in \{0,1\}} f(r_1, \dots, r_{i-1}, x_i, b_{i+1}, \dots, b_\ell)$$

$\mathcal{V}$  checks if  $f_{i-1}(r_{i-1}) = f_i(0) + f_i(1)$  and sends  $r_i \leftarrow \mathbb{F}$ .

(3) In the  $\ell$ th round,  $\mathcal{P}$  sends the polynomial

$$f_\ell(x_\ell) = f(r_1, r_2, \dots, r_\ell)$$

$\mathcal{V}$  checks if  $f_{\ell-1}(r_{\ell-1}) = f_\ell(0) + f_\ell(1)$ . The verifier samples the final random challenge  $r_\ell \in \mathbb{F}$ . Given oracle access to  $f$ ,  $\mathcal{V}$  accepts if and only if  $f_\ell(r_\ell) = f(r_1, r_2, \dots, r_\ell)$ .

### A.3 DNNs and Gradient Descent

We consider two possible layers for the underlying neural network: (i) dense layers, where the linear operation is matrix multiplication; (ii) convolutional layers, where the linear operation is a convolution. A linear operation is followed by an activation, which can be ReLU and tanh for non-output and Softmax for the output layers. In convolutional layers, the activation is additionally followed by an average or max pooling. Moreover, a convolutional layer might include several input and output channels. For instance, if the input to the layer is an RGB image, it has three channels. In such cases, weights are convoluted with each input layer and added together. Each output channel has its own set of weights. We only consider one input/output channel per layer to simplify the algorithm description. The computation is presented in Algorithm 1. In line 20,  $\text{pad}(\mathbf{W}_{i,\ell}^T)$  zero-pads  $\mathbf{W}_{i,\ell}^T$  to an specific dimension and then rotate it for  $\pi$ . More precisely, for  $\mathbf{W}_{i,\ell}$  of size  $w \times w$  and  $\mathbf{U}_{i,\ell-1}$  of size  $n \times n$ ,  $\text{pad}(\mathbf{W}_{i,\ell}^T)$  is of size  $(2w - n) \times (2w - n)$ . In line 24,  $\bar{\mathbf{G}}_{i,l}$  is the average gradient over the data points in the batch.

## B OUR POGD PROTOCOL

In this section, we present the complete PoGD scheme in Protocol 3. Moreover, we formally prove the security of the protocol. We note that the sumcheck runs on the linear operations in step (4) of phase 2 and step (2) of phase 3 and can be merged together. Parties can combine evaluations  $\tilde{\mathbf{T}}_{i,\ell}(r_{2,\mathbf{T}_{i,\ell}}^{(3)})$  and  $\tilde{\mathbf{T}}_{i,\ell}(r_{3,\mathbf{T}_{i,\ell}}^{(1)})$  and then invoke the sumcheck protocol on the combined evaluation. Similarly, scalings in step (6) of phase 2 and step (2) of phase 3 can be merged together.

**Proof of Theorem 1.** Following the completeness of GKR,  $\text{Sum}_{\text{Mat}}$ , and  $\text{Sum}_{\text{Conv}}$ , and PCS completeness of PoGD is immediate. We first show a negligible bound soundness error with a phase-by-phase analyze. Suppose that the prover sends some  $\mathbf{W}_{i,\ell}^* \neq \mathbf{W}_{i,\ell}$ .

**P1:** For  $i = 1$  to  $L$ , following the Schwartz-Zippel lemma and the soundness of GKR, we have  $\tilde{\mathbf{W}}_{i,\ell}^*(r_{\text{out},\mathbf{W}_{i,\ell}}) \neq \tilde{\mathbf{W}}_{i,\ell}(r_{\text{out},\mathbf{W}_{i,\ell}})$  with high probability unless: (i)  $\tilde{\mathbf{W}}_{i-1,\ell}^*(r_{1,\mathbf{W}_{i-1,\ell}}) \neq \tilde{\mathbf{W}}_{i-1,\ell}(r_{1,\mathbf{W}_{i-1,\ell}})$ , which proceeds to phase 4; (ii)  $\tilde{\mathbf{G}}_{i,\ell}^*(r_{1,\mathbf{G}_{i,\ell}}) \neq \tilde{\mathbf{G}}_{i,\ell}(r_{1,\mathbf{G}_{i,\ell}})$ , which proceeds to phase 2.

### Algorithm 1 Mini-Batch Gradient Descent

**Input:** Weights  $\{\mathbf{W}_{i-1,\ell}\}_{j=1}^L$  and Batch  $\mathbf{B}_{i-1}$

**Output:** Updated weights  $\{\mathbf{W}_{i,\ell}\}_{j=1}^L$

#### Forward pass:

Let  $\mathbf{U}_{i,0}$  be the concatenation of input features included in  $\mathbf{B}_{i-1}$

```

1: for  $\ell = 1$  to  $L$  do
2:   if the  $\ell$ -th layer is dense then
3:      $\mathbf{T}_{i,\ell} \leftarrow \mathbf{W}_{i-1,\ell} \cdot \mathbf{U}_{i,\ell-1}$ 
4:      $\mathbf{U}_{i,\ell} \leftarrow \text{act}_\ell(\mathbf{T}_{i,\ell})$   $\triangleright \text{act}_\ell \in \{\text{ReLU}, \text{tanh}, \text{Softmax}\}$ 
5:   else if the  $\ell$ -th layer is convolution then
6:      $\mathbf{T}_{i,\ell} \leftarrow \mathbf{W}_{i-1,\ell} * \mathbf{U}_{i,\ell-1}$ 
7:      $\mathbf{Q}_{i,\ell} \leftarrow \text{act}_\ell(\mathbf{T}_{i,\ell})$   $\triangleright \text{act}_\ell \in \{\text{ReLU}, \text{tanh}, \text{Softmax}\}$ 
8:      $\mathbf{U}_{i,\ell} \leftarrow \text{pool}_\ell(\mathbf{Q}_{i,\ell})$   $\triangleright \text{pool}_\ell \in \{\text{MaxPool}, \text{AvgPool}\}$ 
9:   end if
10: end for

```

#### Backward pass:

Let  $\mathbf{R}_{i,L+1} = \partial \mathcal{L}(\mathbf{U}_{i,L}, \mathbf{Y}_i) / \partial \mathbf{U}_{i,L}$ , s.t.  $\mathbf{Y}_i$  includes labels in  $\mathbf{B}_{i-1}$

```

11: for  $\ell = L$  to 1 do
12:   if the  $\ell$ -th layer is dense then
13:      $\mathbf{T}'_{i,\ell} = \partial \mathbf{U}_{i,\ell} / \partial \mathbf{T}_{i,\ell} = \partial \text{act}_\ell(\mathbf{T}_{i,\ell}) / \partial \mathbf{T}_{i,\ell}$ 
14:      $\mathbf{G}_{i,\ell} = (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell}) \cdot \mathbf{U}_{i,\ell-1}^T$ 
15:      $\mathbf{R}_{i,\ell} = \mathbf{W}_{i-1,\ell}^T \cdot (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell})$ 
16:   else if the  $\ell$ -th layer is convolution then
17:      $\mathbf{Q}'_{i,\ell} = \partial \mathbf{U}_{i,\ell} / \partial \mathbf{Q}_{i,\ell} = \partial \text{pool}_\ell(\mathbf{Q}_{i,\ell}) / \partial \mathbf{Q}_{i,\ell}$ 
18:      $\mathbf{T}'_{i,\ell} = \partial \mathbf{U}_{i,\ell} / \partial \mathbf{T}_{i,\ell} = \mathbf{Q}'_{i,\ell} \circ (\partial \text{act}_\ell(\mathbf{T}_{i,\ell}) / \partial \mathbf{T}_{i,\ell})$ 
19:      $\mathbf{G}_{i,\ell} = (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell}) * \mathbf{U}_{i,\ell-1}$ 
20:      $\mathbf{R}_{i,\ell} = \text{pad}(\mathbf{W}_{i-1,\ell}^T) * (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell})$   $\triangleright$  If  $\ell \neq 1$ 
21:   end if
22: end for

```

#### Update weights:

```

23: for  $\ell = 1$  to  $L$  do
24:    $\mathbf{W}_{i,\ell} = \mathbf{W}_{i-1,\ell} - \eta \cdot \bar{\mathbf{G}}_{i,\ell}$   $\triangleright \eta$  is the learning rate
25: end for

```

**P2:** Suppose that case (ii) is occurred in phase 1 and there exists at least one  $\ell$  such that  $\tilde{\mathbf{G}}_{i,\ell}^*(r_{1,\mathbf{G}_{i,\ell}}) \neq \tilde{\mathbf{G}}_{i,\ell}(r_{1,\mathbf{G}_{i,\ell}})$ . Following the soundness of GKR,  $\text{Sum}_{\text{Mat}}$ , and  $\text{Sum}_{\text{Conv}}$ , there are three possible cases can occur: (i)  $\tilde{\mathbf{W}}_{i-1,\ell}^*(r_{2,\mathbf{W}_{i-1,\ell}}) \neq \tilde{\mathbf{W}}_{i-1,\ell}(r_{2,\mathbf{W}_{i-1,\ell}})$ ,  $\tilde{\mathbf{AUX}}_{i,\ell}^*(r_{2,\mathbf{AUX}_{i,\ell}}) \neq \tilde{\mathbf{AUX}}_{i,\ell}(r_{2,\mathbf{AUX}_{i,\ell}})$ , or for the data batch  $\tilde{\mathbf{B}}_{i,\ell}^*(r_{2,\mathbf{B}_{i,\ell}}) \neq \tilde{\mathbf{B}}_{i,\ell}(r_{2,\mathbf{B}_{i,\ell}})$ , which proceeds to phase 4; The next case is (ii)  $\tilde{\mathbf{R}}_{i,\ell+1}^*(r_{2,\mathbf{R}_{i,\ell+1}}) \neq \tilde{\mathbf{R}}_{i,\ell+1}(r_{2,\mathbf{R}_{i,\ell+1}})$ , proceeding to the next layer; (iii)  $\tilde{\mathbf{U}}_{i,\ell}^*(r_{2,\mathbf{U}_{i,\ell}}) \neq \tilde{\mathbf{U}}_{i,\ell}(r_{2,\mathbf{U}_{i,\ell}})$ , proceeding to phase 3.

**P3:** Suppose (iii) occurs in phase 2 and there exists at least one  $\ell$  such that  $\tilde{\mathbf{U}}_{i,\ell}^*(r_{2,\mathbf{U}_{i,\ell}}) \neq \tilde{\mathbf{U}}_{i,\ell}(r_{2,\mathbf{U}_{i,\ell}})$ . Following the soundness of GKR and  $\text{Sum}_{\text{Mat}}$ , and  $\text{Sum}_{\text{Conv}}$ , there are two possible cases: (i)  $\tilde{\mathbf{W}}_{i-1,\ell}^*(r_{3,\mathbf{W}_{i,\ell}}) \neq \tilde{\mathbf{W}}_{i-1,\ell}(r_{3,\mathbf{W}_{i,\ell}})$ ,  $\tilde{\mathbf{B}}_{i,\ell}^*(r_{3,\mathbf{B}_{i,\ell}}) \neq \tilde{\mathbf{B}}_{i,\ell}(r_{3,\mathbf{B}_{i,\ell}})$ , or  $\tilde{\mathbf{AUX}}_{i,\ell}^*(r_{3,\mathbf{AUX}_{i,\ell}}) \neq \tilde{\mathbf{AUX}}_{i,\ell}(r_{3,\mathbf{AUX}_{i,\ell}})$ , which proceeds phase 4; (ii)  $\tilde{\mathbf{U}}_{i-1,\ell}^*(r_{3,\mathbf{U}_{i-1,\ell}}) \neq \tilde{\mathbf{U}}_{i-1,\ell}(r_{3,\mathbf{U}_{i-1,\ell}})$ , which proceeds to the next iteration, i.e.,  $\ell - 1$ -th layer.

### Protocol 2. GKR

Let  $\mathbb{F}$  be a finite field and  $C : \mathbb{F}^{s_{in}} \rightarrow \mathbb{F}^{s_{out}}$  be a layered arithmetic circuit of depth  $d$ .  $\mathcal{P}$  wants to convince  $\mathcal{V}$  that  $C(\mathbf{in}) = \mathbf{out}$ , where  $\mathbf{in}$  is the input given by  $\mathcal{V}$  and  $\mathbf{out}$  is the output sent to  $\mathcal{V}$ . Without loss of generality, we pad  $s_{in}$  and  $s_{out}$  to powers of 2. The protocol proceeds as follows.

- (1)  $\mathcal{V}$  picks a random  $z \in \mathbb{F}$  and sends it to  $\mathcal{P}$ . Both parties  $\mathcal{P}$  and  $\mathcal{V}$  compute  $\tilde{V}_0(z)$ , where  $\tilde{V}_0$  is the multilinear extension of array  $\mathbf{out}$ .
- (2)  $\mathcal{P}$  and  $\mathcal{V}$  invoke a sumcheck protocol on

$$\tilde{V}_0(z) = \sum_{x, y \in \{0,1\}^{s_1}} \text{add}_0(z, x, y) \cdot (\tilde{V}_1(x) + \tilde{V}_1(y)) + \text{mult}_0(z, x, y) \cdot \tilde{V}_1(x) \cdot \tilde{V}_1(y)$$

At the end of the protocol,  $\mathcal{V}$  receives two evaluations  $\tilde{V}_1(u^{(1)})$ ,  $\tilde{V}_1(v^{(1)})$ .  $\mathcal{V}$  computes  $\text{add}_0(z, u^{(1)}, v^{(1)})$  and  $\text{mult}_0(z, u^{(1)}, v^{(1)})$  and checks that  $\text{add}_0(z, u^{(1)}, v^{(1)}) (\tilde{V}_1(u^{(1)}) + \tilde{V}_1(v^{(1)})) + \text{mult}_0(z, u^{(1)}, v^{(1)}) \cdot \tilde{V}_1(u^{(1)}) \cdot \tilde{V}_1(v^{(1)})$  equals to the last message of sumcheck.

- (3) For  $i = 1, 2, \dots, d$ :

- $\mathcal{V}$  picks two random  $\alpha^{(i)}$  and  $\beta^{(i)}$  and send them to  $\mathcal{P}$ .
- $\mathcal{P}$  and  $\mathcal{V}$  invoke a sumcheck protocol on

$$\alpha^{(i)} \tilde{V}_i(u^{(i)}) + \beta^{(i)} \tilde{V}_i(v^{(i)}) = \sum_{x, y \in \{0,1\}^{s_{i+1}}} \left( \alpha^{(i)} \text{add}_i(u^{(i)}, x, y) + \beta^{(i)} \text{add}_i(v^{(i)}, x, y) \right) \cdot (\tilde{V}_{i+1}(x) + \tilde{V}_{i+1}(y)) \\ + \left( \alpha^{(i)} \text{mult}_i(u^{(i)}, x, y) + \beta^{(i)} \text{mult}_i(v^{(i)}, x, y) \right) \cdot \tilde{V}_{i+1}(x) \cdot \tilde{V}_{i+1}(y)$$

At the end of the protocol,  $\mathcal{V}$  receives two evaluations  $\tilde{V}_{i+1}(u^{(i+1)})$ ,  $\tilde{V}_{i+1}(v^{(i+1)})$ .  $\mathcal{V}$  computes the following and if it does not equal to the last message of the sumcheck, returns 0.

$$\left( \alpha^{(i)} \text{add}_i(u^{(i)}, u^{(i+1)}, v^{(i+1)}) + \beta^{(i)} \text{add}_i(v^{(i)}, u^{(i+1)}, v^{(i+1)}) \right) \cdot (\tilde{V}_{i+1}(u^{(i+1)}) + \tilde{V}_{i+1}(v^{(i+1)})) \\ + \left( \alpha^{(i)} \text{mult}_i(u^{(i)}, u^{(i+1)}, v^{(i+1)}) + \beta^{(i)} \text{mult}_i(v^{(i)}, u^{(i+1)}, v^{(i+1)}) \right) \cdot \tilde{V}_{i+1}(u^{(i+1)}) \cdot \tilde{V}_{i+1}(v^{(i+1)})$$

- (4) At the input layer  $d$ ,  $\mathcal{V}$  has two claims  $\tilde{V}_d(u^{(d)})$  and  $\tilde{V}_d(v^{(d)})$ . As  $\mathcal{V}$  has access to the input  $\mathbf{in}$ , it is possible to directly compute  $\tilde{V}_d$  and evaluate it at two points  $u^d$  and  $v^d$ . If claims are correct,  $\mathcal{V}$  returns 1, otherwise returns 0.

**P4:** Suppose that the case (i) in phases one, two, or three has caused an error, proceeding to this phase. Then, following the soundness of the evaluation reduction sumcheck and also the Schwartz-Zippel lemma, either  $\tilde{W}_{i-1}^*(r_{in}, w_{i-1}) \neq \tilde{W}_{i-1}(r_{in}, w_{i-1})$ ,  $\tilde{B}_i^*(r_{in}, B_i) \neq \tilde{B}_i(r_{in}, B_i)$ , or  $\tilde{AUX}_i^*(r_{in}, AUX_i) \neq \tilde{AUX}_i(r_{in}, AUX_i)$ . Any of these cases cause  $\mathcal{V}$  to return reject based the soundness of the underlying polynomial commitment scheme PCS.

Knowledge soundness follows immediately from the above. As PCS is knowledge sound, the extractor of the PoGD scheme can invoke the extractor of PCS to get witnesses committed to by the prover  $\tilde{W}_i^*$ ,  $\tilde{W}_{i-1}^*$ ,  $\tilde{B}_{i-1}^*$ , and  $\tilde{AUX}_i^*$ . As PoGD is sound, they are satisfying. Moreover, as GKR, Sum<sub>Mat</sub>, Sum<sub>Conv</sub>, and PCS are zero-knowledge, the simulator of PoGD can invoke the simulators of these building blocks to generate proof messages and commitments indistinguishable from a real PoGD protocol execution.

## C COMMITMENT AGGREGATION SCHEME

In this section, we first provide a formal definition of a polynomial commitment aggregation scheme. We then present our construction in Protocol 4 and sketch the security proof. An aggregation scheme consists of a prover  $\mathcal{P}$  and a verifier algorithm  $\mathcal{V}$  such that on input a set of commitment/evaluation instances  $\{(\sigma_i, x_i, y_i)\}_{i=1}^k$  the protocol ends with an aggregated instance  $(\sigma^*, x^*, y^*)$ . Each instance input  $(\sigma_i, x_i, y_i)$  includes a commitment  $\sigma_i$  to a polynomial  $f_i$ , which is supposed to satisfy  $y_i = f_i(x_i)$ . Similarly,  $\sigma^*$  is a commitment to a polynomial  $f^*$  with evaluation  $y^* = f^*(x^*)$ . As the witness, the prover holds input polynomials  $\{f_i\}$  and their corresponding commitment randomnesses. In the end, the prover gets the output polynomial  $f^*$  and its commitment randomness.

If one can substantiate that the aggregated instance  $(\sigma^*, x^*, y^*)$  is satisfied, we can conclude that input instances  $(\sigma_i, x_i, y_i)$  are also satisfied if the aggregation verifier has returned accept. Furthermore, aggregation messages should not reveal any information about each input polynomial  $f_i$  other than its evaluation points  $(x_i, y_i)$ . We next present the security notions more formally.

**Definition 3. (Zero-knowledge aggregation scheme)** Given a knowledge sound zero-knowledge polynomial commitment PCS with parameters  $\text{pp}$ , an aggregation scheme  $\text{AGG} = (\mathcal{P}, \mathcal{V})$  satisfies:

- **Completeness.** Given satisfying instances  $S = \{(\sigma_i, x_i, y_i)\}_{i=1}^k$  with witness  $W = \{(f_i, r_i)\}_{i=1}^k$ , the following probability is 1.

$$\Pr \left[ b_0 \wedge b_1 = 1 \mid \begin{array}{l} (\sigma^*, x^*, y^*), (f^*, r^*), \pi_0 \leftarrow \mathcal{P}(S, W, \text{pp}); \\ b_0 \leftarrow \mathcal{V}(S, (\sigma^*, x^*, y^*), \pi_0, \text{pp}); \\ (y^*, \pi_1) \leftarrow \text{PCS.Open}(f^*, r^*, x^*, \text{pp}); \\ b_1 \leftarrow \text{PCS.Verify}(\sigma^*, x^*, y^*, \pi_1, \text{pp}); \end{array} \right]$$

- **Knowledge soundness.** For any PPT adversary prover  $\mathcal{A}$ , there exist a PPT extractor algorithm  $\xi^{\mathcal{A}}$  such that

$$\Pr \left[ \begin{array}{l} \{(f_i, r_i)\}_{i=1}^k, (f^*, r^*) \leftarrow \xi^{\mathcal{A}}(\text{pp}) : \\ f^*(x^*) = y^* \wedge f_i(x_i) = y_i \quad \forall i \in [k] \end{array} \right] \\ \geq \Pr \left[ \begin{array}{l} \{(\sigma_i, x_i, y_i)\}_{i=1}^k, (\sigma^*, x^*, y^*), \pi_0, \pi_1 \leftarrow \mathcal{A}(\text{pp}); \\ \mathcal{V}(\{(\sigma_i, x_i, y_i)\}_{i=1}^k, (\sigma^*, x^*, y^*), \pi_0, \text{pp}) = 1; \\ \text{PCS.Verify}(\sigma^*, x^*, y^*, \pi_1, \text{pp}) = 1; \end{array} \right] - \text{negl}(\lambda)$$

- **Zero knowledge.** Given a set of tuples  $F = \{(\sigma_i, x_i, y_i)\}_{i=1}^k$ , where each  $f_i$  is a polynomial with evaluation  $f_i(x_i) = y_i$ , there exists a simulator  $\mathcal{S}$  such that given oracle access to the algorithms of PCS, for every adversary distinguisher  $\mathcal{A}$  the following two ideal world and real world experiments are indistinguishable.

### Protocol 3. Proof of Gradient Descent Protocol

**Parameters:** Let GKR be the Virgo++ [62] implementation of a generic GKR-style zero-knowledge proof system, Sum<sub>Mat</sub> be the zero-knowledge sublinear sumcheck protocol for matrix multiplication [53], and Sum<sub>Conv</sub> be the zero-knowledge sublinear sumcheck protocol for convolution [45]. Furthermore, let PCS = (KeyGen, Commit, Open, Verify) be the Orion [60] polynomial commitment scheme with public parameters pp. Suppose that the model has  $L$  layers.

**Initialization:** On inputs weights  $\mathbf{W}_{i-1} = \{\mathbf{W}_{i-1,\ell}\}_{\ell=1}^L$  and a batch  $\mathbf{B}_{i-1}$ ,  $\mathcal{P}$  runs gradient descent (cf. Algorithm 1) to evaluate outputs  $\mathbf{U}_i = \{\mathbf{U}_{i,\ell}\}_{\ell=1}^L$ , gradients  $\mathbf{G}_i = \{\mathbf{G}_{i,\ell}\}_{\ell=1}^L$  and  $\mathbf{R}_i = \{\mathbf{R}_{i,\ell}\}_{\ell=1}^L$ , weights  $\mathbf{W}_i = \{\mathbf{W}_{i,\ell}\}_{\ell=1}^L$ , and auxiliary inputs  $\mathbf{AUX}_i = \{\mathbf{AUX}_{i,\ell}\}_{\ell=1}^L$ .  $\mathcal{P}$  sends  $\mathcal{V}$  the commitments  $\sigma_{W_i} \leftarrow \text{PCS.Commit}(\widetilde{\mathbf{W}}_i, r_{W_i})$ ,  $\sigma_{W_{i-1}} \leftarrow \text{PCS.Commit}(\widetilde{\mathbf{W}}_{i-1}, r_{W_{i-1}})$ ,  $\sigma_{B_{i-1}} \leftarrow \text{PCS.Commit}(\widetilde{\mathbf{B}}_{i-1}, r_{B_{i-1}})$ , and  $\sigma_{\mathbf{AUX}_i} \leftarrow \text{PCS.Commit}(\widetilde{\mathbf{AUX}}_i, r_{\mathbf{AUX}_i})$ .

**Phase 1: Proof of update.** For  $\ell = 1$  to  $L$  do the following steps.

- (1)  $\mathcal{V}$  samples  $r_{\text{out}, \mathbf{W}_{i,\ell}}$  and request the evaluation  $\widetilde{\mathbf{W}}_{i,\ell}(r_{\text{out}, \mathbf{W}_{i,\ell}})$ .
- (2) Having  $\widetilde{\mathbf{W}}_{i,\ell}(r_{\text{out}, \mathbf{W}_{i,\ell}})$ , parties run GKR on  $\mathbf{W}_{i,\ell} = \mathbf{W}_{i-1,\ell} - \eta \cdot \overline{\mathbf{G}}_{i,\ell}$ , which ends with evaluations  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{1, \mathbf{W}_{i-1,\ell}})$ ,  $\overline{\mathbf{G}}_{i,\ell}(r_{1, \mathbf{G}_{i,\ell}})$ .

**Phase 2: Proof of backward pass.** For  $\ell = 1$  to  $L$  do the following steps.

- (1) Having  $\overline{\mathbf{G}}_{i,\ell}(r_{1, \mathbf{G}_{i,\ell}})$  from phase 1, parties run Sum<sub>Conv</sub> on  $\mathbf{G}_{i,\ell} = (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell}) * \mathbf{U}_{i,\ell-1}$  if the  $\ell$  layer is convolutional or otherwise Sum<sub>Mat</sub> on  $\mathbf{G}_{i,\ell} = (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell}) \cdot \mathbf{U}_{i,\ell-1}^T$  if the layer is dense. At the end of the run, parties receive evaluations  $\widetilde{\mathbf{T}}'_{i,\ell+1}(r_{2, \mathbf{T}'_{i,\ell+1}}^{(1)})$ ,  $\widetilde{\mathbf{U}}_{i,\ell}(r_{2, \mathbf{U}_{i,\ell}}^{(1)})$ , and  $\widetilde{\mathbf{R}}_{i,\ell+1}(r_{2, \mathbf{R}_{i,\ell+1}}^{(1)})$ .
- (2) For  $\ell \neq 1$ , having  $\widetilde{\mathbf{R}}_{i,\ell}(r_{2, \mathbf{R}_{i,\ell}})$  from the previous iteration, parties run Sum<sub>Conv</sub> on  $\mathbf{R}_{i,\ell} = \text{pad}(\mathbf{W}_{i-1,\ell}^T) * (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell})$  if the layer is convolutional or Sum<sub>Mat</sub> on  $\mathbf{R}_{i,\ell} = \mathbf{W}_{i-1,\ell}^T \cdot (\mathbf{R}_{i,\ell+1} \circ \mathbf{T}'_{i,\ell})$  if the layer is dense, ending with evaluations  $\widetilde{\mathbf{R}}_{i,\ell+1}(r_{2, \mathbf{R}_{i,\ell+1}}^{(2)})$ ,  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{2, \mathbf{W}_{i-1,\ell}}^{(2)})$ , and  $\widetilde{\mathbf{T}}'_{i,\ell+1}(r_{2, \mathbf{T}'_{i,\ell+1}}^{(2)})$ .
- (3) Parties combine evaluations  $\widetilde{\mathbf{T}}'_{i,\ell+1}(r_{2, \mathbf{R}_{i,\ell+1}}^{(1)})$  and  $\widetilde{\mathbf{T}}'_{i,\ell}(r_{2, \mathbf{T}'_{i,\ell}}^{(2)})$ , by running an evaluation reduction sumcheck, i.e., Equation (5). Having the combination, parties run GKR on  $\mathbf{T}'_{i,\ell} = \mathbf{Q}'_{i,\ell} \circ (\partial \text{act}_\ell(\mathbf{T}_{i,\ell}) / \partial \mathbf{T}_{i,\ell})$  followed by another run on  $\mathbf{Q}'_{i,\ell} = \partial \text{pool}_\ell(\text{act}_\ell(\mathbf{T}_{i,\ell})) / \partial \text{act}_\ell(\mathbf{T}_{i,\ell})$  if the layer is convolutional or otherwise  $\mathbf{T}'_{i,\ell} = \partial \mathbf{U}_{i,\ell} / \partial \mathbf{T}_{i,\ell} = \partial \text{act}_\ell(\mathbf{T}_{i,\ell}) / \partial \mathbf{T}_{i,\ell}$  if the layer is dense. At the end, parties receive evaluations  $\widetilde{\mathbf{T}}_{i,\ell}(r_{2, \mathbf{T}_{i,\ell}}^{(3)})$  and  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{2, \mathbf{AUX}_{i,\ell}}^{(3)})$ .
- (4) Having  $\widetilde{\mathbf{T}}_{i,\ell}(r_{2, \mathbf{T}_{i,\ell}}^{(3)})$ , parties run Sum<sub>Conv</sub> on  $\mathbf{T}_{i,\ell} = \mathbf{W}_{i-1,\ell} * \mathbf{U}_{i,\ell-1}$  if the layer is convolutional or Sum<sub>Mat</sub> on  $\mathbf{T}_{i,\ell} = \mathbf{W}_{i-1,\ell} \cdot \mathbf{U}_{i,\ell-1}$  if the layer is dense. At the end, parties received evaluations  $\widetilde{\mathbf{U}}_{i,\ell-1}(r_{2, \mathbf{U}_{i,\ell-1}}^{(4)})$  and  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{2, \mathbf{W}_{i-1,\ell}}^{(4)})$ . For  $\ell = 1$ ,  $\widetilde{\mathbf{U}}_{i,0}(r_{2, \mathbf{U}_{i,0}}^{(4)})$  is viewed as  $\widetilde{\mathbf{B}}_{i-1}(r_{2, \mathbf{B}_{i-1}}^{(4)})$ .
- (5) For  $\ell = L$ , parties combine  $\widetilde{\mathbf{R}}_{i,L+1}(r_{2, \mathbf{R}_{i,L+1}}^{(1)})$  and  $\widetilde{\mathbf{R}}_{i,L+1}(r_{2, \mathbf{R}_{i,L+1}}^{(2)})$  by an evaluation reduction sumcheck. Having the combinations, parties run GKR on  $\mathbf{R}_{i,L+1} = \partial \mathcal{L}(\mathbf{U}_{i,L}, \mathbf{Y}_i) / \partial \mathbf{U}_{i,L}$ . At the end, parties receive evaluations  $\widetilde{\mathbf{B}}_{i-1}(r_{2, \mathbf{B}_{i-1}}^{(5)})$ ,  $\widetilde{\mathbf{U}}_{i,L}(r_{2, \mathbf{U}_{i,L}}^{(5)})$ , and  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{2, \mathbf{AUX}_{i,\ell}}^{(5)})$ .
- (6) Consistency checks and scalings are proved by running sumchecks on Equation (3) and Equation (4), which ends with several evaluations of auxiliary inputs  $\widetilde{\mathbf{AUX}}_{i,\ell}$ . The evaluations are combined together into a single evaluation  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{2, \mathbf{AUX}_{i,\ell}}^{(6)})$ .
- (7) For  $\ell = L$ , evaluations of  $\widetilde{\mathbf{W}}_{i-1,\ell}$ ,  $\widetilde{\mathbf{B}}_{i-1}$ ,  $\widetilde{\mathbf{U}}_{i,\ell}$ ,  $\widetilde{\mathbf{AUX}}_{i,\ell}$ , from steps 1-5, are combined into  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{2, \mathbf{W}_{i-1,\ell}})$ ,  $\widetilde{\mathbf{U}}_{i,\ell}(r_{2, \mathbf{U}_{i,\ell}})$ ,  $\widetilde{\mathbf{B}}_{i-1}(r_{2, \mathbf{B}_{i-1}})$ ,  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{2, \mathbf{AUX}_{i,\ell}})$ .

**Phase 3: Proof of forward pass.** For  $\ell = L$  to 1 do the following steps.

- (1) Let  $\widetilde{\mathbf{U}}_{i,L}(r_{3, \mathbf{U}_{i,L}}) = \widetilde{\mathbf{U}}_{i,L}(r_{2, \mathbf{U}_{i,L}})$ . Parties combine  $\widetilde{\mathbf{U}}_{i,\ell}(r_{3, \mathbf{U}_{i,\ell}})$  received from the previous iteration (if  $\ell \neq L$ ) and  $\widetilde{\mathbf{U}}_{i,\ell-1}(r_{2, \mathbf{U}_{i,\ell-1}})$  from phase 2. Having the combined evaluation, parties run GKR on  $\mathbf{U}_{i,\ell} = \text{pool}_\ell(\mathbf{Q}_{i,\ell})$  followed by another run on  $\mathbf{Q}_{i,\ell} = \text{act}_\ell(\mathbf{T}_{i,\ell})$  if the layer is convolutional or on  $\mathbf{U}_{i,\ell} = \text{act}_\ell(\mathbf{T}_{i,\ell})$  if the layer is dense, ending with  $\widetilde{\mathbf{T}}_{i,\ell}(r_{3, \mathbf{T}_{i,\ell}}^{(1)})$  and  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{3, \mathbf{AUX}_{i,\ell}}^{(1)})$ .
- (2) Having  $\widetilde{\mathbf{T}}_{i,\ell}(r_{3, \mathbf{T}_{i,\ell}}^{(1)})$ , parties run Sum<sub>Conv</sub> on  $\mathbf{T}_{i,\ell} = \mathbf{W}_{i-1,\ell} * \mathbf{U}_{i,\ell-1}$  or Sum<sub>Mat</sub> on  $\mathbf{T}_{i,\ell} = \mathbf{W}_{i-1,\ell} \mathbf{U}_{i,\ell-1}$ . At the end, parties receive evaluations  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{3, \mathbf{W}_{i-1,\ell}})$  and  $\widetilde{\mathbf{U}}_{i,\ell-1}(r_{3, \mathbf{U}_{i,\ell-1}})$ . If  $i = 1$ ,  $\widetilde{\mathbf{U}}_{i,0}(r_{3, \mathbf{U}_{i,0}})$  is viewed as  $\widetilde{\mathbf{B}}_{i-1,\ell}(r_{3, \mathbf{B}_{i-1,\ell}})$ . For consistency checks and scalings, parties run sumchecks on Equation (3) and Equation (4), ending with several evaluations combined together and with  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{3, \mathbf{AUX}_{i,\ell}}^{(1)})$  into  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{3, \mathbf{AUX}_{i,\ell}})$ .

**Phase 4: Proof of evaluation reduction.**

- (1) Parties combine received evaluations of form  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{1, \mathbf{W}_{i-1,\ell}})$ ,  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{2, \mathbf{W}_{i-1,\ell}})$ , and  $\widetilde{\mathbf{W}}_{i-1,\ell}(r_{3, \mathbf{W}_{i-1,\ell}})$  into a single evaluation  $\widetilde{\mathbf{W}}_{i-1}(r_{\text{in}, \mathbf{W}_{i-1}})$ , evaluations of form  $\widetilde{\mathbf{B}}_{i-1,\ell}(r_{2, \mathbf{B}_{i-1,\ell}})$ ,  $\widetilde{\mathbf{B}}_{i-1,\ell}(r_{3, \mathbf{B}_{i-1,\ell}})$  into a single evaluation  $\widetilde{\mathbf{B}}_{i-1}(r_{\text{in}, \mathbf{B}_{i-1}})$  and evaluations of form  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{1, \mathbf{AUX}_{i,\ell}})$ ,  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{2, \mathbf{AUX}_{i,\ell}})$ , and  $\widetilde{\mathbf{AUX}}_{i,\ell}(r_{3, \mathbf{AUX}_{i,\ell}})$  into a single evaluation  $\widetilde{\mathbf{AUX}}_i(r_{\text{in}, \mathbf{AUX}_i})$ . The final combined evaluations can be verified by invoking the evaluation opening algorithm of the commitment scheme  $\text{PCS.Open}(\widetilde{\mathbf{W}}_i, r_{W_i}, r_{\text{in}, \mathbf{AUX}_i}, \text{pp})$ ,  $\text{PCS.Open}(\widetilde{\mathbf{W}}_{i-1}, r_{W_{i-1}}, r_{\text{in}, \mathbf{W}_{i-1}}, \text{pp})$ ,  $\text{PCS.Open}(\widetilde{\mathbf{B}}_{i-1}, r_{B_{i-1}}, r_{\text{in}, \mathbf{B}_{i-1}}, \text{pp})$ , and  $\text{PCS.Open}(\widetilde{\mathbf{AUX}}_i, r_{\mathbf{AUX}_i}, r_{\text{in}, \mathbf{AUX}_i}, \text{pp})$  or applying the aggregation in the recursive setting (cf. Section 6).

#### Protocol 4. Polynomial Commitment Aggregation

**Parameters:** Let PCS = (KeyGen, Commit, Open, Verify) be the Orion [60] polynomial commitment scheme with public parameters pp. Let  $f_1, f_2, \dots, f_k$  be  $\ell$ -variate polynomials with variable-degree  $d$ , each has an evaluation  $y_i = f_i(x_i)$  and committed to by  $\sigma_i$  with randomness  $r_i$ .  $\mathbf{M}_i$  denotes coefficients of  $f_i$  parsed as a square matrix and  $C_{i,1}, C_{i,2}$  denote the associated encoded matrices generated during the Orion commitment procedure (cf. Section 5.2).

##### Phase 1: Evaluation Aggregation

- (1)  $\mathcal{P}$  and  $\mathcal{V}$  interpolate a polynomial  $L : \mathbb{F} \rightarrow \mathbb{F}^\ell$  such that  $L(i) = x_i$ .
- (2)  $\forall i = 1$  to  $k$ ,  $\mathcal{P}$  samples random univariate polynomials  $h_{i,1}, \dots, h_{i,\ell}$ , each of degree  $d\ell$ , sets  $h_i(z_1, \dots, z_\ell) \leftarrow h_{i,1}(z_1) + \dots + h_{i,\ell}(z_\ell)$ , and sends  $v_i \leftarrow h_i(x_i)$  and  $\sigma'_i \leftarrow \text{PCS.Commit}(h_i, r'_i, \text{pp})$ .  $\mathbf{M}'_i$  denotes the coefficient matrix of  $h_i$  and  $C'_{i,1}, C'_{i,2}$  denote the associated encoded matrices.
- (3)  $\forall i = 1$  to  $k$ ,  $\mathcal{V}$  responds with a random challenge  $\beta_i$ .
- (4)  $\forall i = 1$  to  $k$ ,  $\mathcal{P}$  sends the polynomial  $g_i \leftarrow (f_i + \beta_i h_i) \circ L$ .
- (5)  $\forall i = 1$  to  $k$ ,  $\mathcal{V}$  checks whether  $g_i(i) = y_i + \beta_i v_i$  and responds with random challenges  $\alpha_i$  and  $r$ .
- (6) Parties compute  $g \leftarrow \sum_{i=1}^k \alpha_i g_i$  and return  $x^* = L(r)$ , and  $y^* = g(r)$  as output evaluation points.
- (7)  $\mathcal{P}$  computes the output witness polynomial  $f^* = \sum_{i=1}^k \alpha_i f_i + \alpha_i \beta_i h_i$ .

##### Phase 2: Commitment Aggregation

- (1)  $\mathcal{P}$  sends  $\sigma^* \leftarrow \text{PCS.Commit}(f^*, r^*, \text{pp})$ .  $\mathbf{M}^*$  denotes the coefficient matrix of  $f^*$  and  $C_1^*, C_2^*$  denote the associated encoded matrices.
- (2)  $\mathcal{V}$  randomly samples and sends a subset of indices  $I$  of size  $t \in \Theta(\lambda)$ .
- (3)  $\forall i = 1$  to  $k$ ,  $\forall \text{idx} \in I$ ,  $\mathcal{P}$  sends  $s^* = C_2^*[\text{idx}]$ ,  $s_i = C_{i,2}[\text{idx}]$ , and  $s'_i = C'_{i,2}[\text{idx}]$  with their corresponding Merkle proofs.
- (4)  $\forall i = 1$  to  $k$ ,  $\forall \text{idx} \in I$ ,  $\mathcal{V}$  whether  $s^* = \sum_{i=1}^k \alpha_i s_i + \alpha_i \beta_i s'_i$  and the received Merkle proofs are valid.
- (5) The verifier returns  $(\sigma^*, x^*, y^*)$  as the aggregate instances, the prover returns  $(\sigma^*, x^*, y^*), (f^*, r^*)$  as the aggregated instance/witness tuple.

$\text{Real}_{\mathcal{A},F}(1^\lambda, \text{pp})$

- (1)  $\forall i \in \{1, 2, \dots, k\}$ ,  $\sigma_i \leftarrow \text{PCS.Commit}(f_i, r_i, \text{pp})$
- (2) Let  $S = \{(\sigma_i, x_i, y_i)\}_{i=1}^k$  and  $W = \{(f_i, r_i)\}_{i=1}^k$ .
- (3)  $(\sigma^*, x^*, y^*), (f^*, r^*), \pi_0 \leftarrow \mathcal{P}(S, W, \text{pp})$
- (4)  $(y^*, \pi_1) \leftarrow \text{PCS.Open}(f^*, r^*, x^*, \text{pp})$ ;
- (5)  $b \leftarrow \mathcal{A}(S, (\sigma^*, x^*, y^*), \pi_0, \pi_1, \text{pp})$
- (6) Return  $b$

$\text{Ideal}_{\mathcal{A},S}(1^\lambda, \text{pp})$

- (1)  $S = \{(\sigma_i, x_i, y_i)\}_{i=1}^k, (\sigma^*, x^*, y^*), \pi_0, \pi_1 \leftarrow \mathcal{S}(1^\lambda, \text{pp})$
- (2)  $b \leftarrow \mathcal{A}(S, (\sigma^*, x^*, y^*), \pi_0, \pi_1, \text{pp})$
- (3) Return  $b$

In particular, for any  $\mathcal{A}$ , we have

$$|\Pr[\text{Real}_{\mathcal{A},F}(1^\lambda, \text{pp}) = 1] - \Pr[\text{Ideal}_{\mathcal{A},S}(1^\lambda, \text{pp}) = 1]| \leq \text{negl}(\lambda)$$

**Proof of Theorem 2.** Completeness is followed immediately. The soundness proof proceeds in two steps. First, we must prove the soundness with respect to commitment aggregation, i.e., if we can prove the well-formedness of  $\sigma^*$ , the verifier is convinced that each  $\sigma_i$  is well-formed, and the linear combination relation holds. Second, we should prove the soundness with respect to evaluation aggregation, i.e., if a prover can substantiate  $(\sigma^*, x^*, y^*)$  is satisfied, the verifier is convinced that  $\forall i : (\sigma_i, x_i, y_i)$  is satisfied.

Consider the case of commitment aggregation. Let  $\delta$  be the distance of the underlying code. First, suppose that the linear combination relation holds. Assume there exists some  $\sigma_i$  or  $\sigma'_i$ , which is not well-formed, i.e., there exists a  $c$ th column of  $C_{i,2}$  or  $C'_{i,2}$ , which is not a valid codeword. Then, over the choice of random

challenges  $\alpha_i$  and  $\beta_i$ , the  $c$ th column of  $C_2$  is not a valid codeword with overwhelming probability [32, Claim 1], which yields  $\sigma^*$  to be not well-formed. Similarly, if there exist an  $r$ th row of  $C_{i,1}$  or  $C'_{i,1}$ , which is not code-word, it would cause the  $r$ th row of  $C_1$  to be an invalid code-word, yielding the verifier to return reject during the well-formedness check of  $\sigma^*$ . We should then show that the verifier returns reject if the linear combination does not hold. Suppose that this case occurs, and we have  $\mathbf{M}^* = \mathbf{M} + \sum_i \alpha_i \mathbf{M}_i + \alpha_i \beta_i \mathbf{M}'_i$ , where  $\mathbf{M}$  is a non-zero matrix, i.e., there exist an index  $(\text{idx}_r, \text{idx}_c)$  such that  $\mathbf{M}[\text{idx}_r, \text{idx}_c] \neq 0$ . If they apply the row-wise encoding, we have  $C_1^* = C_1 + \sum_i \alpha_i C_{i,1} + \alpha_i \beta_i C'_{i,1}$ , where the row  $C_1[\text{idx}_r, :]$  has at least  $\delta$  non-zero entries. Then, if we apply the column-wise encoding, we have  $C_2^* = C_2 + \sum_i \alpha_i C_{i,2} + \alpha_i \beta_i C'_{i,2}$ , where  $C_2$  includes  $\delta^2$  non-zero entries. In other words, the two-step application of the encoding scheme propagates the inconsistency through the matrix. Let  $\gamma = \frac{\delta}{p}$ , be the relative distance of the code, where  $p$  is the codeword size. By opening  $\Theta(\lambda)$  entries, all opened values are zero with a negligible probability of at most  $(1 - \gamma^2)^{\Theta(\lambda)}$ .

Then, for the case of evaluation aggregation, suppose that we have  $\exists i : f_i(x_i) \neq y_i$ . If the prover sends a correct  $g_i = (f_i + \beta_i h_i) \circ L$ , then with overwhelming probability over the choice of the  $\beta_i$ , we should have  $y_i + \beta_i v_i \neq g_i(i)$ , and the verifier aborts. Therefore, the prover must send some  $g_i \neq (f_i + \beta_i h_i) \circ L$  to pass the verification. However, with overwhelming probability over the choice of the  $\alpha_i$ ,

$$g = \sum_i \alpha_i g_i \neq \sum_i \alpha_i (f_i + \beta_i h_i) \circ L = f \circ L$$

where  $f = \sum_{i=1}^k \alpha_i f_i + \alpha_i \beta_i h_i$ . The Schwartz-Zippel lemma then implies  $(f \circ L)(r) \neq g(r)$  except with negligible probability.

Knowledge soundness follows immediately from the above. Our aggregation scheme requires the prover to give the verifier oracle access to matrices  $C_2^*$ ,  $C_{i,2}$ , and  $C'_{i,2}$  by their Merkle root. It is well known that for such Merkle-based oracle accesses, one can build a straight-line extractor to extract the witness committed to by its Merkle root [9, 48]. Once matrices are extracted, we can recover coefficients by running the decoding procedure first to columns and then to rows of each matrix.

Zero knowledge holds since each polynomial  $h_{i,j}$  is uniformly sampled from the set of polynomials in  $\mathbb{F}[X]$  of degree  $d\ell$ . The simulator  $\mathcal{S}$  first samples random polynomials  $\{h_{i,j}^*\}_{i=1,j=1}^{k,\ell}$  and sets  $h_i^*(z_1, \dots, z_\ell) = h_{i,1}^*(z_1) + \dots + h_{i,\ell}^*(z_\ell)$ . Then,  $\mathcal{S}$  commits to each  $h_i^*$  and sends  $\sigma_i^*$  and the evaluation  $v_i^* = h_i^*(x)$  to  $\mathcal{A}$ . As  $h_{i,j}^*$  in the ideal world and  $h_{i,j}$  in the real world both are sampled uniformly from the same distribution, they are indistinguishable, and hence, polynomials  $h_i^*$  and  $h_i$  and the evaluations  $v_i$  and  $v_i^*$  are indistinguishable. Moreover, as PCS is zero-knowledge  $\sigma_i'$  and  $\sigma_i'^*$  are indistinguishable.  $\mathcal{S}$  then samples input polynomials  $\{f_i^*\}_{i=1}^k$  such that  $f_i^*(x_i) = y_i$  and sends  $\sigma_i^*$ , a commitment to  $f_i^*$  to the verifier, which is indistinguishable from  $\sigma_i$ .  $\mathcal{S}$  also sends polynomials  $g_i^* = (f_i^* + \beta_i h_i^*) \circ L$ . We note that each  $f_i \circ L$  is a univariate polynomial from the set  $\mathbb{F}[X]$  of degree at most  $dk\ell$ , and  $h_i \circ L$  is also a univariate polynomial over the field with the same degree. Therefore, each coefficient of the polynomial  $f_i \circ L$ , and hence  $g_i$ , is masked by a coefficient of  $h_i \circ L$ . Then, as masks  $h_i^*$  and  $h_i$  are indistinguishable, thus  $g_i^*$  and  $g_i$  are also indistinguishable.

## D RECURSIVE SUMCHECK FRAMEWORK

In this section, we present our recursive sumcheck framework in Protocol 5 and sketch its security proof.

**Proof of Theorem 3.** For completeness, consider a satisfying tuple  $(i-1, z_0, z_{i-1}, \pi_{i-1})$ , i.e., we have  $\mathcal{V}(i-1, z_0, z_{i-1}, \pi_{i-1}, \text{pp}) = 1$ . Given  $\omega_{i-1}$ , we must show that the updated instance  $(i, z_0, z_i, \pi_i)$ , returned by the prover is also satisfying.  $\mathcal{P}$  executes the  $i$ th iteration and generates  $\pi_i = (\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ . First, as  $\mathcal{V}(i-1, z_0, z_{i-1}, \pi_{i-1}, \text{pp}) = 1$ , the verification steps, and in turn, the entire computation of  $\mathcal{F}_A$  is completed successfully.  $\pi_{i,G}$  is generated by GKR.  $\mathcal{P}$  and the sumcheck protocol, and  $\pi_{i,A}$  is generated by AGG.  $\mathcal{P}$ . If the iteration is a final iteration  $\pi_{i,E}$  is generated by PCS.  $\mathcal{P}$  as an opening to  $\text{agg}_{i+1}$ , otherwise it is returned by AGG.  $\mathcal{P}$  as a witness to  $\text{agg}_{i+1}$ . Moreover,  $\text{wp}_{i+1}$  is the combined predicate evaluation by applying evaluation reduction sumcheck. Following the completeness of the GKR scheme, the sumcheck protocol, AGG proof messages  $\pi_{i,G}$  and  $\pi_{i,A}$  are satisfying. The completeness of AGG or PCS ensures that  $\pi_{i,E}$  is also satisfying. The correctness of  $\text{wp}_{i+1}$  follows the completeness of the sumcheck protocol. Therefore, the  $i$ th proof  $\pi_i$  is satisfying.

For knowledge soundness, let  $\xi_{\text{AGG}}$  be the extractor of AGG. The adversary  $\mathcal{A}$  returns a satisfying tuple  $(i, z_0, z_i, \pi_i)$ , i.e., we have  $\mathcal{V}(i, z_0, z_i, \pi_i, \text{pp}) = 1$ . We must construct an efficient extractor  $\xi$  that can extract witnesses that correctly satisfy  $\mathcal{F}$  for the entire computation. We show inductively that for all  $j = i-1$  to 1,  $\xi$  can construct  $\xi^{(j)}$  that outputs  $z_j$  and  $\omega_j$ . For the base case, i.e.,  $j = i-1$ , recall that  $\pi_i = (\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ . As  $\pi_{i,A}$  is a valid proof for  $\text{agg}_{i+1}$ , by running  $\xi_{\text{AGG}}$ ,  $\xi^{(i-1)}$  can extract the

vector  $\mathbf{u}_i = i \|z_0\|z_{i-1}\|\omega_{i-1}\|\pi_{i-1}$ , which is committed to by  $\sigma_i$ . Therefore,  $z_{i-1}$  and  $\omega_{i-1}$  are extracted. Soundness of GKR, PCS, AGG, and the protocol sumcheck implies that the extracted inputs  $\mathbf{u}_i$ , so  $z_{i-1}$  and  $\omega_{i-1}$  are satisfying and  $z_i = \mathcal{F}(z_{i-1}, \omega_{i-1})$ . For the case  $j < i-1$ ,  $\xi^{(j)}$  runs  $\xi^{(j+1)}$  to extract  $\mathbf{u}_{j+1} = j+1 \|z_0\|z_j\|\omega_j\|\pi_j$ . As  $\mathbf{u}_{j+1}$  is satisfying, proof messages in  $\pi_j$ , including  $\pi_{j,A}$  are valid. By running  $\xi_{\text{AGG}}$ ,  $\xi^{(j)}$  can then extract  $\mathbf{u}_j$  and get  $z_{j-1}$  and  $\omega_{j-1}$ . This completes the recursive construction of  $\xi$  from  $\xi_{\text{AGG}}$ .

We prove zero-knowledge with respect to the *final* verifier. In particular, we construct a simulator  $\mathcal{S}$  that can generate proof messages commitments indistinguishable from the real proofs and commitments. In particular,  $\mathcal{S}$  samples a random initial input  $z_0$  and auxiliary inputs  $\{\omega_j\}_{j=0}^{i-1}$ . With access to the  $\mathcal{F}$ ,  $\mathcal{S}$  executes the IVC prover algorithm  $i$  iterations. At the end,  $\mathcal{S}$  returns the final tuple  $(i, z_0, z_i, \pi_i)$ . If the instantiations of GKR, PCS, AGG, and the sumcheck protocol satisfy zero-knowledge, proof messages and commitments included in  $\pi_i$  reveals no additional information about the private auxiliary inputs  $\{\omega_j\}_{j=0}^{i-1}$ .

## E PROOF OF TRAINING

In this section, we first define zero-knowledge proof-of-training. We then present the KAIZEN construction in Protocol 6.

### E.1 Definition

A zero-knowledge proof-of-training must satisfy completeness, knowledge soundness, and zero-knowledge defined formally as below. Let  $C$  be the training circuit defined in Section 7.2.

**Definition 4. (Zero-Knowledge Proof-of-Training)** A tuple of (KeyGen, DataCom, WeiCom, BatchOpen, Prove, Verify) is a zkPoT if the following hold. Let  $\text{pp} \leftarrow \text{KeyGen}(1^\lambda)$ .

- **Completeness.** Given a dataset  $\mathcal{D}$ ,  $\rho_{\mathcal{D}} \leftarrow \text{DataCom}(\mathcal{D}, \text{pp}, r_{\rho_{\mathcal{D}}})$ ,  $\mathbf{W}_0$  and  $\mathbf{W}_{i-1}$  with commitments  $\sigma_{\mathbf{W}_0} \leftarrow \text{WeiCom}(\mathbf{W}_0, r_{\sigma_{\mathbf{W}_0}}, \text{pp})$  and  $\sigma_{\mathbf{W}_{i-1}} \leftarrow \text{WeiCom}(\mathbf{W}_{i-1}, r_{\sigma_{\mathbf{W}_{i-1}}}, \text{pp})$ , respectively, the and opening  $\mathbf{B}_{i-1}, \mathbf{p}_{\mathbf{B}_{i-1}} \leftarrow \text{BatchOpen}(\mathcal{D}, r_{\rho_{\mathcal{D}}}, i, \text{pp})$ , and a satisfying proof  $\pi_{i-1}$  the following equals to 1. Let  $t_i = (i, \rho_{\mathcal{D}}, \sigma_{\mathbf{W}_0}, \sigma_{\mathbf{W}_{i-1}})$ .

$$\Pr \left[ b = 1 \mid \begin{array}{l} \mathbf{W}_i, \pi_i \leftarrow \text{Prove}(t_i, \mathbf{W}_{i-1}, \mathbf{B}_{i-1}, \mathbf{p}_{\mathbf{B}_{i-1}}, \pi_{i-1}, \text{pp}); \\ \mathbf{W}_i \leftarrow C(\mathbf{W}_{i-1}, \mathbf{B}_{i-1}, (i, \rho_{\mathcal{D}}, \mathbf{p}_{\mathbf{B}_{i-1}})); \\ \sigma_{\mathbf{W}_i} \leftarrow \text{WeiCom}(\mathbf{W}_i, r_{\sigma_{\mathbf{W}_i}}, \text{pp}); \\ b \leftarrow \text{Verify}(i, \rho_{\mathcal{D}}, \sigma_{\mathbf{W}_0}, \sigma_{\mathbf{W}_i}, \pi_i, \text{pp}); \end{array} \right]$$

- **Knowledge Soundness.** For any adversary  $\mathcal{A}$ , there exists an extractor  $\xi^{\mathcal{A}}$  such that the following holds.

$$\Pr \left[ \begin{array}{l} \{\mathbf{B}_k, \mathbf{p}_{\mathbf{B}_k}, \mathbf{W}_k\}_{k=0}^i \leftarrow \xi^{\mathcal{A}}(\text{pp}) : \\ \forall k \in \{1, \dots, k\} : \mathbf{W}_k = C(\mathbf{W}_{k-1}, \mathbf{B}_{k-1}, k, \rho_{\mathcal{D}}, \mathbf{p}_{\mathbf{B}_{k-1}}); \\ \sigma_{\mathbf{W}_i} = \text{WeiCom}(\mathbf{W}_i, r_{\sigma_{\mathbf{W}_i}}, \text{pp}); \sigma_{\mathbf{W}_0} = \text{WeiCom}(\mathbf{W}_0, r_{\sigma_{\mathbf{W}_0}}, \text{pp}); \end{array} \right] \geq \Pr \left[ \begin{array}{l} i, \rho_{\mathcal{D}}, \sigma_{\mathbf{W}_0}, \sigma_{\mathbf{W}_i}, \pi_i \leftarrow \mathcal{A}(\text{pp}) : \\ \text{Verify}(i, \rho_{\mathcal{D}}, \sigma_{\mathbf{W}_0}, \sigma_{\mathbf{W}_i}, \pi_i, \text{pp}) = 1 \end{array} \right] - \text{negl}(\lambda)$$

- **Zero knowledge.** There exists a simulator  $\mathcal{S}$ , having oracle access to  $C$ , such that for every adversary  $\mathcal{A}$ , iteration counter  $i$ , dataset  $\mathcal{D}$  with a commitment  $\rho_{\mathcal{D}}$ , initial weights  $\mathbf{W}_0$ , batches of data points and openings  $\{\mathbf{B}_k, \mathbf{p}_k\}_{k=0}^i$  satisfying the circuit  $C$ , the two following experiments Real, Ideal are indistinguishable.

### Protocol 5. Recursive Sumcheck Framework

**Parameters:** Let  $\text{PCS} = (\text{KeyGen}, \text{Commit}, \text{Open}, \text{Verify})$  be the Orion [60] commitment with public parameters  $\text{pp}$ ,  $\text{AGG} = (\text{Prove}, \text{Verify})$  the non-interactive aggregation scheme (cf. Protocol 4), and  $\text{GKR} = (\text{Prove}, \text{Verify})$  be a generic GKR-style non-interactive zero-knowledge proof. Hashes are instantiated with  $\text{MiMC-p/p}$  [4]. Let  $\mathcal{F}$  be the iteration function,  $z_0$  an initial input,  $\pi_0$  a trivially satisfying proof, and  $\omega_{i-1}$  the  $i$ th auxiliary input.

$z_i, \pi_i \leftarrow \mathcal{P}(i, z_0, z_{i-1}, \omega_{i-1}, \pi_{i-1}, \text{pp}) :$

- (1) Parse  $\pi_{i-1}$  as  $(\pi_{i-1,G}, \pi_{i-1,A}, \pi_{i-1,E}, \text{agg}_{i-1}, \text{wp}_{i-1}, \text{agg}_i, \text{wp}_i)$ , where  $\pi_{i-1,G}$  is sumcheck and  $\pi_{i-1,A}$  is aggregation messages,  $\text{agg}_{i-1}$  and  $\text{agg}_i$  are the last and updated input aggregation instances,  $\text{wp}_{i-1}$  and  $\text{wp}_i$  are the last and updated wiring predicate aggregations, and  $\pi_{i-1,E}$  is a witness to  $\text{agg}_i$ .
- (2) Compute the augmented iteration function  $i, z_0, z_i, \text{agg}_i, \text{wp}_i \leftarrow \mathcal{F}_A(i, z_0, z_{i-1}, \omega_{i-1}, \pi_{i-1} \setminus \{\pi_{i-1,E}\})$  as follows:
  - (2.1) Evaluate Fiat-Shamir and Merkle hashes required for verification of sumcheck proof messages  $\pi_{i-1,G}$  and aggregation proof messages  $\pi_{i-1,A}$ .
  - (2.2) On input  $i-1, z_0, z_{i-1}, \text{agg}_{i-1}, \text{wp}_{i-1}$ , run  $\text{GKR.V}$  and the sumcheck protocol verifier to verify sumcheck messages  $\pi_{i-1,G}$ . The verification must end with an input commitment/evaluation instance  $(\sigma_{i-1}, x_{i-1}, y_{i-1})$  and the updated predicate aggregation  $\text{wp}_i$ .
  - (2.3) On input  $\text{agg}_{i-1}, \text{agg}_i, (\sigma_{i-1}, x_{i-1}, y_{i-1})$ , run  $\text{AGG.V}$  to verify aggregation messages  $\pi_{i-1,A}$ .
  - (2.4) If verification steps (2.1)-(2.3) passes, compute  $z_i \leftarrow \mathcal{F}(z_{i-1}, \omega_{i-1})$ , and return  $i, z_0, z_i, \text{agg}_i, \text{wp}_i$ .
- (3) Evaluate a commitment  $\sigma_i$  to the inputs of the running  $\mathcal{F}_A$ , i.e., given  $\mathbf{u}_i = i || z_0 || z_{i-1} || \omega_{i-1} || \pi_{i-1} \setminus \{\pi_{i-1,E}\}$ , we have  $\sigma_i \leftarrow \text{PCS.Commit}(\tilde{\mathbf{u}}_i, r_i, \text{pp})$ .
- (4) Generate sumcheck messages for the running  $\mathcal{F}_A$  execution by applying hash sumcheck, i.e., sumchecks on Equation (6), to hash evaluations and GKR.P to other components. Combine wiring predicate evaluations into single one  $\text{wp}_{i+1}$  and also the evaluations of  $\tilde{\mathbf{u}}_i$  into  $(x_i, y_i)$  by applying evaluation reduction sumchecks, i.e., running the sumcheck protocol on Equation (5). Include sumcheck messages as well as  $\text{wp}_{i+1}$  and  $(\sigma_i, x_i, y_i)$  in  $\pi_{i,G}$ .
- (5) On input instances  $\text{agg}_i, (\sigma_i, x_i, y_i)$  and witnesses  $\pi_{i-1,E}, (\tilde{\mathbf{u}}_i, r_i)$ , run  $\text{AGG.P}$ , which ends with  $\text{agg}_{i+1}$ . Include aggregation messages in  $\pi_{i,A}$ .
- (6) If the running iteration is the final iteration, run  $\text{PCS.Open}$  and include evaluation opening proofs otherwise the witnesses for  $\text{agg}_{i+1}$  in  $\pi_{i,E}$ .
- (7) Let  $\pi_i \leftarrow (\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ , and return  $z_i, \pi_i$ .

$\{0, 1\} \leftarrow \mathcal{V}(i, z_0, z_i, \pi_i, \text{pp}) :$

- (1) Parse  $\pi_i$  as  $(\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ .
- (2) Evaluate Fiat-Shamir and Merkle hashes required for verification of sumcheck proof messages  $\pi_{i,G}$  and aggregation proof messages  $\pi_{i,A}$ .
- (3) On input  $i, z_0, z_i, \text{agg}_i, \text{wp}_i$  run  $\text{GKR.V}$  and the sumcheck protocol verifier on messages  $\pi_{i,G}$ , which ends with  $(\sigma_i, x_i, y_i)$  and  $\text{wp}_{i+1}$ . Having access to wiring predicates, verify  $\text{wp}_{i+1}$  directly. On input  $\text{agg}_{i+1}, \text{agg}_i, (\sigma_i, x_i, y_i)$  run  $\text{AGG.V}$  to verify messages  $\pi_{i,A}$ . If it is the final iteration, run  $\text{PCS.Verify}(\text{agg}_{i+1}, \pi_{i,E}, \text{pp})$ . Otherwise, check whether  $\pi_{i,E}$  is a valid opening witness to  $\text{agg}_{i+1}$ .

Let  $\text{INP}$  be a set, including the dataset, initial weights, batches with their satisfying openings generated by  $\text{BatchOpen}$ .

$\text{Real}_{\mathcal{A}, \text{INP}}(i, \text{pp})$

- (1)  $\rho_{\mathcal{D}} \leftarrow \text{DataCom}(\mathcal{D}, \text{pp}, \rho_{\mathcal{D}})$
- (2)  $\sigma_{W_0} \leftarrow \text{WeiCom}(W_0, r_{\sigma_{W_0}}, \text{pp})$
- (3)  $\forall k$ , given  $t_k = (k, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_{k-1}})$ :
  - (3.1)  $W_k, \pi_k \leftarrow \text{Prove}(t_k, W_{k-1}, B_{k-1}, p_{B_{k-1}}, \pi_{k-1}, \text{pp})$
  - (3.2)  $\sigma_{W_k} \leftarrow \text{WeiCom}(W_k, r_{\sigma_{W_k}}, \text{pp})$
- (4)  $b \leftarrow \mathcal{A}(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}, \pi_i, \text{pp})$
- (5) Return  $b$

$\text{Ideal}_{\mathcal{A}, \mathcal{S}}(i, \text{pp})$

- (1)  $\rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}, \pi_i \leftarrow \mathcal{S}(\text{pp})$ , given oracle access to  $\mathcal{C}$
- (2)  $b \leftarrow \mathcal{A}(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}, \pi_i, \text{pp})$
- (3) Return  $b$

In particular, for any instance  $t_i$  and adversary  $\mathcal{A}$  we have

$$|\Pr[\text{Real}_{\mathcal{A}, \text{INP}}(i, \text{pp}) = 1] - \Pr[\text{Ideal}_{\mathcal{A}, \mathcal{S}}(i, \text{pp}) = 1]| \leq \text{negl}(\lambda)$$

## E.2 Construction

We present the  $\text{KAIZEN}$  construction in Protocol 6. As  $\text{KAIZEN}$  is basically a concrete instantiation of Protocol 5, we extend the recursion protocol to describe the prover and verifier of  $\text{KAIZEN}$ .

Also, we sketch the security proof, which follows the same approach as the security proof, presented in Appendix D.

**Proof of Theorem 4.** For completeness, consider a satisfying instance/proof such that  $\text{Verify}(i-1, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_{i-1}}, \pi_{i-1}, \text{pp}) = 1$ . Given weights  $W_{i-1}$  and batch  $B_{i-1} \subseteq \mathcal{D}$  with opening  $p_{B_{i-1}}$ , we show that the updated instance,  $(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i})$  and proof  $\pi_i$  should be satisfying.  $\mathcal{P}$  executes the  $i$ th iteration and generates  $\pi_i = (\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ . Note that as we have  $\text{Verify}(i-1, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_{i-1}}, \pi_{i-1}, \text{pp}) = 1$ , the verification steps of  $\mathcal{F}_A$  are completed successfully. Furthermore, the completeness of the Merkle tree ensures that the computation  $\mathcal{C}$  is also completed successfully. Then,  $\pi_{i,G}$  is generated by  $\text{GKR.P}$ ,  $\text{PoGD.P}$ , and the sumcheck protocol, and  $\pi_{i,A}$  is generated by  $\text{AGG.P}$ . If the iteration is a final iteration  $\pi_{i,E}$  is generated by  $\text{PCS.P}$ , otherwise, by  $\text{AGG.P}$  as an opening or witness to  $\text{agg}_{i+1}$ . Moreover,  $\text{wp}_{i+1}$  is the combined predicate evaluation by applying evaluation reduction sumcheck. Following the completeness of  $\text{GKR}$ ,  $\text{PoGD}$ , the sumcheck protocol, and  $\text{AGG}$ , the messages included in  $\pi_{i,G}$  and  $\pi_{i,A}$  are satisfying. The completeness of  $\text{AGG}$  or  $\text{PCS}$  in the case of the final iteration, ensures that  $\pi_{i,E}$  is satisfying. Moreover, the correctness of  $\text{wp}_{i+1}$  follows the completeness of the sumcheck protocol.

For knowledge soundness, let  $\xi_{\text{AGG}}$  be the extractor of  $\text{AGG}$ . The adversary  $\mathcal{A}$  returns a satisfying tuple  $(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}, \pi_i)$ , i.e., we have  $\text{Verify}(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}, \pi_i, \text{pp}) = 1$ . We must construct

### Protocol 6. KAIZEN Zero-Knowledge Proof-of-Training

**Parameters:** Let PCS = (KeyGen, Commit, Open, Verify) be the Orion [60] commitment with public parameters pp, AGG = (Prove, Verify) the non-interactive aggregation scheme (cf. Protocol 4), PoGD = (Prove, Verify) the non-interactive zero-knowledge proof of gradient descent (cf. Protocol 3), and GKR be a non-interactive zero-knowledge generic GKR-style sumcheck-based proof system [62]. Let  $\mathcal{D}$  be the dataset,  $W_0$  be an initial weights,  $C$  be the training iteration circuit as defined in Section 7.2. Let  $\pi_0$  be a trivially satisfying proof.

**Basic procedures:** KeyGen samples public parameters of the commitment scheme by running PCS.KeyGen, DataCom generates a commitment  $\rho_{\mathcal{D}}$  to the dataset  $\mathcal{D}$  by returning its Merkle root. WeiCom commits the multilinear extension of weights by apply PCS.Commit, and BatchOpen generates an opening proof for a batch  $B$ , consistent with the permutation hard-coded in  $C$ , by returning Merkle paths of each point included in the batch.

$W_i, \pi_i \leftarrow \text{Prove}(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_{i-1}}, W_{i-1}, B_{i-1}, p_{B_{i-1}}, \pi_{i-1}, pp) :$

- (1) Parse  $\pi_{i-1}$  as  $(\pi_{i-1,G}, \pi_{i-1,A}, \pi_{i-1,E}, \text{agg}_{i-1}, \text{wp}_{i-1}, \text{agg}_i, \text{wp}_i)$ , where  $\pi_{i-1,G}$  is sumcheck and  $\pi_{i-1,A}$  is aggregation messages,  $\text{agg}_{i-1}$  and  $\text{agg}_i$  are the last and updated input aggregation instances,  $\text{wp}_{i-1}$  and  $\text{wp}_i$  are the last and updated wiring predicate aggregations, and  $\pi_{i-1,E}$  is a witness to  $\text{agg}_i$ .
- (2) Compute the augmented iteration function  $i, \rho_{\mathcal{D}}, \sigma_{W_0}, W_i, \text{agg}_i, \text{wp}_i \leftarrow \mathcal{F}_A(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_{i-1}}, W_{i-1}, B_{i-1}, p_{B_{i-1}}, \pi_{i-1} \setminus \{\pi_{i-1,E}\})$  as follows.
  - (a) Evaluate Fiat-Shamir and Merkle hashes required for verification of sumcheck proof messages  $\pi_{i-1,G}$  and aggregation proof messages  $\pi_{i-1,A}$ .
  - (b) On input  $i-1, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_{i-1}}, \text{agg}_{i-1}, \text{wp}_{i-1}$ , verify sumcheck messages  $\pi_{i-1,G}$  by running GKR.V, PoGD.V, and the sumcheck verifier. The verification ends with several commitment-evaluation instances of form  $\{(\sigma_{i-1,j}, x_{i-1,j}, y_{i-1,j})\}_j$  and the updated predicate aggregation  $\text{wp}_i$ , where each  $\sigma_{i-1,j}$  is a commitment to either  $\bar{W}_{i-2}, \bar{B}_{i-2}$ , gradient descent auxiliary inputs  $\text{AUX}_{i-1}$ , or other inputs, which we denote them by  $u_{i-1} = i-1 \|\rho_{\mathcal{D}}\| \|\sigma_{W_0}\| \|\sigma_{W_{i-2}}\| \|\text{p}_{B_{i-2}}\| \|\pi_{i-2}$ . Moreover,  $\sigma_{i-1,j}$  can be a commitment to  $\bar{W}_{i-1}$ , which should be the same as  $\sigma_{W_{i-1}}$ .
  - (c) On input  $\text{agg}_{i-1}, \{(\sigma_{i-1,j}, x_{i-1,j}, y_{i-1,j})\}_j$ , verify aggregation messages  $\pi_{i-1,A}$  by running AGG.V. The verification ends with the updated aggregation instances  $\text{agg}_i$ ; note that each type of polynomials, i.e., either  $\bar{W}, \bar{B}, \text{AUX}$ , or  $\bar{u}$ , are aggregated independently.
  - (d) Compute  $W_i = C(W_{i-1}, B_{i-1}, (i, \rho_{\mathcal{D}}, p_{B_{i-1}}))$ , and return  $i, \rho_{\mathcal{D}}, \sigma_{W_0}, W_i, \text{agg}_i, \text{wp}_i$ .
- (3) Generate commitments  $\sigma_{i,0} \leftarrow \text{PCS.Commit}(\bar{W}_i, r_{\sigma_{W_i}}, pp)$ ,  $\sigma_{i,1} \leftarrow \text{PCS.Commit}(\bar{B}_{i-1}, r_{\sigma_{B_{i-1}}}, pp)$ ,  $\sigma_{i,2} \leftarrow \text{PCS.Commit}(\text{AUX}_i, r_{\sigma_{\text{AUX}_i}}, pp)$ , where  $\text{AUX}_i$  is the auxiliary inputs of the running gradient descent iteration, and  $\sigma_{i,3} \leftarrow \text{PCS.Commit}(\bar{u}_i, r_{\sigma_{\bar{u}_i}}, pp)$ . Let  $\sigma_{i,4} = \sigma_{W_{i-1}}$ .
- (4) Generate sumcheck messages for the running  $\mathcal{F}_A$  execution by applying hash sumcheck, i.e., sumchecks on Equation (6), to hash evaluations, PoGD.P to the gradient descent computation, and GKR.P to any other components. Moreover, combine received wiring predicate evaluations into  $\text{wp}_{i+1}$  by applying evaluation reduction sumcheck, i.e., sumcheck on Equation (5). Similarly, combine all evaluations of  $\bar{u}_i$ . Include sumcheck messages, commitments generated in step (3), and evaluations required for sumcheck messages in  $\pi_{i,G}$ .
- (5) On input  $\text{agg}_i, \{(\sigma_i, x_i, y_i)\}_j$  with witnesses  $\bar{W}_i, \bar{W}_{i-1}, \bar{B}_{i-1}, \text{AUX}_i$ , or  $\bar{u}_i$ , and  $\pi_{i-1,E}$ , run AGG.P to get  $\text{agg}_{i+1}$ . Include messages in  $\pi_{i,A}$ .
- (6) If the running iteration is the final iteration, run PCS.Open and include evaluation opening proofs otherwise the witnesses for  $\text{agg}_{i+1}$  in  $\pi_{i,E}$ .
- (7) Return  $W_i, \pi_i \leftarrow (\pi_{i-1,G}, \pi_{i-1,A}, \pi_{i-1,E}, \text{agg}_{i-1}, \text{wp}_{i-1}, \text{agg}_i, \text{wp}_i)$ .

$\{0, 1\} \leftarrow \text{Verify}(i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}, \pi_i, pp) :$

- (1) Parse  $\pi_i$  as  $(\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ .
- (2) Evaluate Fiat-Shamir and Merkle hashes required for verification of sumcheck proof messages  $\pi_{i,G}$  and aggregation proof messages  $\pi_{i,A}$ .
- (3) On input  $i, \rho_{\mathcal{D}}, \sigma_{W_0}, \sigma_{W_i}$ , run GKR.V, PoGD.V, and sumcheck verifier to verify sumcheck messages  $\pi_{i,G}$ , yielding  $\{(\sigma_{i-1,j}, x_{i-1,j}, y_{i-1,j})\}_j$  and  $\text{wp}_{i+1}$  at the end. Having access to wiring predicates, verify  $\text{wp}_{i+1}$  directly. Moreover, run AGG.V on input  $\text{agg}_i, \{(\sigma_{i-1,j}, x_{i-1,j}, y_{i-1,j})\}_j$  to verify aggregation messages  $\pi_{i,A}$ , yielding  $\text{agg}_{i+1}$  at the end. If the  $i$ th iteration is the final iteration, run PCS.Verify( $\text{agg}_{i+1}, \pi_{i,E}$ ). Otherwise, check whether  $\pi_{i,E}$  is a valid opening witness to  $\text{agg}_{i+1}$ .

an extractor  $\xi$  that can extract witnesses that satisfy  $C$  iterations. Inductively, for  $j = i-1$  to 1,  $\xi$  runs  $\xi^{(j)}$  outputting the  $j$ th witness. Consider  $j = i-1$ , s.t.  $\pi_i = (\pi_{i,G}, \pi_{i,A}, \pi_{i,E}, \text{agg}_i, \text{wp}_i, \text{agg}_{i+1}, \text{wp}_{i+1})$ . As  $\pi_{i,E}$  is a valid opening to  $\text{agg}_{i+1}$ , by running  $\xi_{\text{AGG}}$ ,  $\xi^{(i-1)}$  can extract weights  $W_i, W_{i-1}$ , batch  $B_{i-1}$ , and  $p_{B_{i-1}}$  such that the iteration circuit satisfies  $W_i = C(W_{i-1}, B_{i-1}, (j, \rho, p_{B_{i-1}}))$  and also we have  $\sigma_{W_i} = \text{WeiCom}(W_0, r_{\sigma_{W_i}}, pp)$ . This follows the soundness of GKR, PCS, AGG, and the sumcheck protocol as well as binding of the Merkle tree. Moreover,  $\xi^{(i-1)}$  can extract the inputs of the iteration, i.e.,  $u_i = i \|\rho_{\mathcal{D}}\| \|\sigma_{W_0}\| \|\sigma_{W_{i-1}}\| \|\text{p}_{B_{i-1}}\| \|\pi_{i-1}$ . For the case  $j < i-1$ ,  $\xi^{(j)}$  runs  $\xi^{(j+1)}$  to extract  $u_{j+1} = j+1$ , including proof messages in  $\pi_j$ , so in turn, including  $\pi_{j,E}$ . By running  $\xi_{\text{AGG}}$ , we can then extract  $W_i, W_{i-1}$ , batch  $B_{i-1}, p_{B_{i-1}}$ , and  $u_j$  as before. Finally, by running  $\xi^{(1)}$ , which in turn runs previous extractors  $\xi^{(2)}, \dots, \xi^{(i-1)}$ ,  $\xi$  can

output all witnesses for the entire computation, i.e., intermediate model weights and all data points used in training iterations.

We prove zero knowledge with respect to the *final* verifier.  $\mathcal{S}$  samples a random dataset  $\mathcal{D}^*$  and commits to it by a Merkle root. As Merkle commitments are hiding, the root returned by  $\mathcal{S}$  is indistinguishable from  $\rho$ . Also,  $\mathcal{S}$  samples initial weights  $W_0^*$ . With access to the  $C$ ,  $\mathcal{S}$  executes the prover algorithm on weights and satisfying batches for  $i$  iterations. At the end  $\mathcal{S}$  returns the final proof messages and commitments to  $\sigma_{W_0^*}$  and  $\sigma_{W_i^*}$ . As PCS, AGG are zero-knowledge, commitments  $\sigma_{W_0^*}, \sigma_{W_i^*}$ , and other commitments included in the proof, e.g., aggregated commitment, with their openings are indistinguishable from real commitments. Moreover, as AGG, PoGD, GKR, and the sumcheck protocol are zero-knowledge,  $\pi_i^*$  is also indistinguishable for a real proof.